

Learning Mixtures of Arbitrary Distributions over Large Discrete Domains

Yuval Rabani*

Leonard J. Schulman†

Chaitanya Swamy‡

Abstract

We give an algorithm for learning a mixture of *unstructured* distributions. This problem arises in various unsupervised learning scenarios, for example in learning *topic models* from a corpus of documents spanning several topics. We show how to learn the constituents (the topic distributions and the mixture weights) of a mixture of k (constant) arbitrary distributions over a large discrete domain $[n] = \{1, 2, \dots, n\}$, using $O(n \text{ polylog } n)$ samples.

This task is information-theoretically impossible for $k > 1$ under the usual sampling process from a mixture distribution. However, there are situations (such as the above-mentioned topic model case) in which each sample point consists of several observations from the same mixture constituent. This number of observations, which we call the “*sampling aperture*”, is a crucial parameter of the problem. We show that efficient learning is possible exactly at the information-theoretically least-possible aperture of $2k - 1$. (Independent work by others places certain restrictions on the model, which enables learning with smaller aperture, albeit using, in general, a significantly larger sample size.)

A sequence of tools contribute to the algorithm, such as concentration results for random matrices, dimension reduction, moment estimations, and sensitivity analysis.

1 Introduction

We give an algorithm for learning a mixture of *unstructured* distributions. More specifically, we consider the problem of learning a mixture of k arbitrary distributions over a large finite domain $[n] = \{1, 2, \dots, n\}$. This finds applications in various unsupervised learning scenarios including *collaborative filtering* [26], and learning *topic models* from a corpus of documents spanning several topics [36, 11], which we will use as our prototypical motivating example. Our goal is to learn the probabilistic model that is hypothesized to generate the observed data. In particular, we learn the constituents of the mixture, i.e., the k distributions defining the topics, and their weights in the mixture.

It is information-theoretically impossible to reconstruct the mixture model from single-snapshot samples (e.g., single-word documents). Thus, our work relies on multi-snapshot samples. To illustrate, in the (pure documents) topic model introduced in [36], each document is consists of a *bag of words* generated by selecting a topic with probability proportional to its mixture weight and then taking independent samples from this topic’s distribution (over words); so n is the size of the vocabulary and k is the number of topics. Notice that typically n will be quite large, and

*The Rachel and Selim Benin School of Computer Science and Engineering and the Center of Excellence on Algorithms, The Hebrew University of Jerusalem, Jerusalem 91904, Israel. Email: yraani@cs.huji.ac.il.

†Caltech, Pasadena, CA 91125, USA. Supported in part by NSF CCF-1038578, NSF CCF-0515342, NSA H98230-06-1-0074, and NSF ITR CCR-0326554. Email: schulman@caltech.edu.

‡Dept. of Combinatorics and Optimization, Univ. Waterloo, Waterloo, ON N2L 3G1, Canada. Supported in part by NSERC grant 32760-06. Email: cswamy@math.uwaterloo.ca.

significantly larger than k . Also, clearly, if very long documents are available, the problem becomes trivial, as each document already provides a very good sample for the distribution of its topic. Thus, it is desirable to keep the dependence of the sample size on n as low as possible, while at the same time minimize what we call the *aperture*, which is the number of snapshots per sample point (i.e., words per document). These parameters govern both the applicability of an algorithm and its computational complexity.

Our results. Let $p^1, \dots, p^k \in \Delta^{n-1}$ denote the k -mixture constituents, where Δ^{n-1} is the $(n-1)$ -simplex, and w_1, \dots, w_k denote the mixture weights. Our algorithm uses

$$O\left(\frac{k^3 n \ln n}{\epsilon^6}\right) + O\left(\frac{k^2 n \ln^6 n \ln\left(\frac{k}{\epsilon}\right)}{\epsilon^4}\right) + O\left(\frac{k}{\epsilon}\right)^{O(k^2)} \quad (1)$$

documents (i.e., samples) and reconstructs with high probability (see Theorem 4.1) each mixture constituent up to ℓ_1 -error ϵ , and each mixture weight up to additive error ϵ . We make no assumptions on the constituents. The asymptotic notation hides factors that are polynomial in $w_{\min} := \min_t w_t$ and the “width” of the mixture (which intuitively measures the minimum variation distance between any two constituents). The three terms in (1) correspond to the requirements for the number of 1-, 2-, and $(2k-1)$ -snapshots respectively. So we need aperture $2k-1$ only for a small part of the sample. (Clearly, longer documents can be split into pieces that can be used as independent samples.)

To put our bounds in perspective, notice importantly that we recover the mixture constituents within ℓ_1 -distance ϵ . With fixed aperture (independent of n), a sample size of $\Omega(n)$ is *necessary* to recover even the expectation of the mixture distribution with constant ℓ_1 -error. On the other hand, aperture $\Omega((n+k^2)\log nk)$ is sufficient for algorithmically trivial recovery of the model with constant ℓ_∞ error using few samples. Restricting the aperture to $2k-2$ makes recovery impossible to arbitrary accuracy (without additional assumptions): there could be two far-apart k -mixtures that generate exactly the same sample distribution. Thus, we obtain near-optimal dependence on n and optimal aperture.

Our work provides new insights into the widely-studied problem of learning topic models and nicely complements the recent interesting work of [5, 4, 3]. This body of work recovers the constituents (under certain assumptions) up to ℓ_2 or ℓ_∞ error, using a sample size that is $\text{poly}(k)$ and independent of (or sublinear in) n . However, if we seek to achieve ℓ_1 -error ϵ , there are inputs for which their sample size (although $\text{poly}(n, k)$) is $\Omega(n^3)$ (or worse, again ignoring dependence on w_{\min} and “width”; see Appendix A). This is a significantly poorer dependence on n compared to our near-linear dependence (so our bounds are better when n is quite large but k is small, e.g., a constant). Observe that with $\Omega(n^3)$ samples, the entire distribution on 3-word documents can be estimated fairly accurately; the challenge in [5, 4, 3] is therefore to recover the model from this relatively noiseless data. In contrast, a major challenge that we face to achieve ℓ_1 -reconstruction with $O(n \text{ polylog } n)$ samples is to ensure that the error remains bounded despite the presence of very noisy data due to the small sample size, and we need to develop suitable machinery to achieve this. (An interesting research direction would be to combine the various approaches to obtain a $O(n \cdot \text{poly}(k, \ln n))$ sample size.)

We now give a rough sketch of our algorithm (Algorithm 1 in Section 3) and the ideas behind its analysis (Section 4). Let $P = (p^1, \dots, p^k)$, $r = \sum_t w_t p^t$ be the expectation of the mixture, and $k' = \text{rank}(p^1 - r, \dots, p^k - r)$. Our algorithm reduces the problem to the problem of learning *one-dimensional mixtures*. We choose k' random lines that are close to the affine hull, $\text{aff}(P)$, of P and “project” the mixture on to these k' lines. We learn each projected mixture, which is a

one-dimensional mixture-learning problem, and then combine the inferred projections on these k' lines to obtain k points that are close to $\text{aff}(P)$. Finally, we project these k' points on to Δ^{n-1} to obtain k distributions over $[n]$, which we argue are close (in ℓ_1 -distance) to p^1, \dots, p^k .

Various difficulties arise in implementing this plan. We first learn a good approximation to $\text{aff}(P)$ using spectral techniques and only 2-snapshots. We use ideas similar to [32, 7, 31], but our challenge is to show that the covariance matrix $A = \sum_t w_t (p^t - r)(p^t - r)^\dagger$ can be well-approximated by the empirical covariance matrix with only $O(n \ln^6 n)$ 2-snapshots. A random orthonormal basis of the learned affine space then supplies the k' lines on which we project our mixture. Of course, we do not know P , so “projecting” on to a basis vector b actually means that we project snapshots from P on to b by mapping item i to b_i . For this to be meaningful, we need to ensure that if the mixture constituents are far apart in variation distance then their projections $(b^\dagger p^t)_{t \in [k]}$ are also well separated relative to the spread of the support $\{b_1, \dots, b_n\}$ of the one-dimensional distribution. In general, we can only claim a relative separation of $\Theta(\frac{1}{\sqrt{n}})$ (since $\min_{t \neq t'} \|p^t - p^{t'}\|_2$ may be $\Theta(\frac{1}{\sqrt{n}})$). We avoid this via a careful balancing act: we prove (Lemma 4.3) that the ℓ_∞ norm of unit vectors in $\text{aff}(P)$ is $O(\frac{1}{\sqrt{n}})$, and argue that this isotropy property suffices since b is close to $\text{aff}(P)$.

Finally, a key ingredient of our algorithm (see Section 5) is to show how to learn the real projections $(b^\dagger p^t)_{t \in [k]}$ from the projected snapshots. This is technically the most difficult step and the one that requires aperture $2k - 1$ (the smallest aperture at which this is information-theoretically possible). We show that the projected snapshots on b yield empirical moments of a related distribution and use this to learn the projections and the mixture weights via a method of moments (see, e.g., [23, 22, 28, 10, 33, 4]). One technical difficulty is that variation distance in Δ^{n-1} translates to transportation distance [39] in the one-dimensional projection. We use a combination of convex programming and numerical-analysis techniques to learn the projections from the empirical “directional” moments. In the process, we establish some novel properties about the *moment curve*—an object that plays a central role in convex and polyhedral geometry [8]—that may be of independent interest.

Related work. The past decade has witnessed tremendous progress in the theory of learning statistical mixture models. The most striking example is that of learning mixtures of high dimensional Gaussians. Starting with Dasgupta’s groundbreaking paper [19], a long sequence of improvements [20, 6, 38, 29, 1, 22, 13] culminated in the recent results [28, 10, 33] that essentially resolve the problem in its general form. In this vein, other highly structured mixture models, such as mixtures of discrete product distributions [30, 24, 17, 23, 14, 16] and similar models [17, 9, 34, 29, 18, 15, 21], have been studied intensively. One important difference between this line of work and ours is that the structure of those mixtures enables learning using single-snapshot samples, whereas this is impossible in our case. Another interesting difference between our setting and the work on structured models (and this is typical of most results on PAC-style learning) is that the amount of information in each sample point is roughly in the same ballpark as the information needed to describe the model. In our setting, the amount of information in each sample point is exponentially sparser than the information needed to describe the model to good accuracy. Thus, the topic-modeling problem motivates the natural question of inference from sparse samples. This issue is also encountered in collaborative filtering; see [31] for some related theoretical problems.

Recently, we learnt about an independent line of inquiry into very much the same question as ours [5, 4, 3]¹. All three papers make certain assumptions about the mixture constituents, which

¹An earlier stage of this work, including the case of $k = 2$ topics as well as some other results that are not subsumed by this paper, dates to 2007. The last version of that phase has been posted since May 2008 at [37]. The extension

makes it possible to learn the mixture constituents with constant aperture. In comparison with our work, their $\text{poly}(n, k)$ sample size (for ℓ_1 -error) is attractive in terms of k but has a worse dependence on n ($\Omega(n^3)$).

The assumptions in [5, 4, 3] impose some limitations on the applicability of their algorithms. To understand this, it is illuminating to consider the case where all the p^t s lie on a line-segment in Δ^{n-1} . This poses no problems for our algorithm, and we recover the p^t s along with their mixture weights. However, as we show below, the algorithms in [5, 4, 3] all fail to reconstruct this mixture. Anandkumar et al. [4] solve the same problem that we consider, under the assumption that P (viewed as an $n \times k$ matrix) has rank k . This assumption is clearly violated in this setting, rendering their algorithm inapplicable. The other two papers [5, 3] deal with the more general setting where each document is generated by a combination of topics [36, 25]: first a convex combination $\lambda \in \Delta^{k-1}$ is sampled from a mixture distribution \mathcal{T} on Δ^{k-1} , then the document is generated by sampling words from the distribution $\sum_{t=1}^k \lambda_t p^t$. The goal is to learn the topic distributions and the mixture distribution. (The problem we consider is the special case where \mathcal{T} places weight w^t on the t -th vertex of Δ^{k-1} .) [5] posits a ρ -separability assumption on the topics, wherein each topic p^t has a unique *anchor word* i such that $p_i^t \geq \rho$ and $p_i^{t'} = 0$ for every $t' \neq t$, whereas [3] weakens this to the requirement that the p^t s be linearly independent. Both papers show how to learn the model when \mathcal{T} is the Dirichlet distribution (which gives the latent Dirichlet model [12]); the paper [5] obtains results for other mixture distributions as well.

In order to apply these algorithms, we can view the input as being specified by two topics, x and y , which are the end points of the line segment; \mathcal{T} then places weight w_t on the convex combination $(\lambda_t, 1 - \lambda_t)^\dagger$, where $p^t = \lambda_t x + (1 - \lambda_t)y$. This \mathcal{T} is far from the Dirichlet distribution, so [4] does not apply here. Suppose that x and y satisfy the ρ -separability condition. (Note that ρ may only be $O(\frac{1}{n})$, even if x and y have *disjoint* supports.) We can then apply the algorithm of Arora et al. [5]. But this *does not* recover \mathcal{T} ; it returns the topic-topic correlation matrix $E_{\mathcal{T}}[\lambda\lambda^\dagger]$, which does not reconstruct the mixture (w, P) .

This limitation should not be surprising since [5] uses constant aperture. Indeed, [5] notes that it is impossible to reconstruct \mathcal{T} with arbitrary accuracy (with any constant aperture) even if one knows the topics x and y . In this context, we remark that our earlier work [37] uses the approach presented in this paper and solves the problem for documents that are *arbitrary* mixtures of two topics, yielding a crisp statement about the tradeoff between the sampling aperture and the accuracy with which \mathcal{T} can be learnt.

Finally, it is also pertinent to compare the topic modeling problem with the problem of learning a mixture of product distributions (e.g., [23]). Multi-snapshot samples can be thought of as single-snapshot samples from the power distribution on $[n]^K$, where K is the aperture. The product distribution literature typically deals with samples spaces that are the product of many small cardinality components, whereas the topic modeling problem deals with samples spaces that are the product of few large cardinality components.

2 Preliminaries

2.1 Mixture sources, snapshots, and projections

Let $[n]$ denote $\{1, 2, \dots, n\}$, and Δ^{n-1} denote the $(n-1)$ -simplex $\{x \in \mathbb{R}_+^n : \sum_i x_i = 1\}$. A k -mixture source (w, P) on $[n]$ consists of k mixture constituents $P = (p^1, p^2, \dots, p^k)$, where p^t has support $[n]$ for all $t \in [k]$, along with the corresponding mixture weights $w = (w_1, \dots, w_k) \in \Delta^{k-1}$.

to arbitrary k is from this year.

An m -snapshot from (w, P) is obtained by choosing $t \in [k]$ according to the distribution w , and then choosing $i \in [n]$ m times independently according to the distribution p^t . The probability distribution on m -snapshots is thus a mixture of k power distributions on the product space $[n]^m$.

We also consider mixture sources whose constituents are distributions on \mathbb{R} . A k -mixture source (w, P) on \mathbb{R} consists of k mixture constituents $P = (p^1, p^2, \dots, p^k)$, where each p^t is a probability distribution on \mathbb{R} , along with corresponding mixture weights $w = (w_1, \dots, w_k) \in \Delta^{k-1}$.

Given a distribution p on $[n]$ and a vector $x \in \mathbb{R}^n$, we define the projection of p on x , denoted $\pi_x(p)$, to be the discrete distribution on \mathbb{R} that assigns probability mass $\sum_{i: x_i = \beta} p_i$ to $\beta \in \mathbb{R}$. (Thus, $\pi_x(p)$ has support $\{x_1, \dots, x_n\}$ and $\mathbb{E}[\pi_x(p)] = x^\dagger p$.) Given a k -mixture source (w, P) on $[n]$, we define the projected k -mixture source $(w, \pi_x(P))$ on \mathbb{R} to be the k -mixture source on \mathbb{R} given by $(w, (\pi_x(p^1), \dots, \pi_x(p^k)))$.

We also denote by $(w, \mathbb{E}[\pi_x(P)])$ the distribution that assigns probability mass w_t to $\mathbb{E}[\pi_x(p^t)] = x^\dagger p^t$ for all $t \in [k]$. This is an example of what we call a k -spike distribution, which is a distribution on \mathbb{R} that assigns positive probability mass to k points in \mathbb{R} .

2.2 Transportation distance for mixtures

Let $(w, (p^1, \dots, p^k))$ be a k -mixture source on $[n]$, and $(\tilde{w}, (\tilde{p}^1, \dots, \tilde{p}^\ell))$ be an ℓ -mixture source on $[n]$. The *transportation distance* (with respect to the total variation distance $\frac{1}{2}\|x - y\|_1$ on measures on Δ^{n-1}) between these two mixture sources, denoted by $\text{Tran}(w, P; \tilde{w}, \tilde{P})$, is the optimum value of the following linear program (LP).

$$\min \sum_{i=1}^k \sum_{j=1}^\ell x_{ij} \cdot \frac{1}{2} \|p^i - \tilde{p}^j\|_1 \quad \text{s.t.} \quad \sum_{j=1}^\ell x_{ij} = w_i \quad \forall i \in [k], \quad \sum_{i=1}^k x_{ij} = \tilde{w}_j \quad \forall j \in [\ell], \quad x \geq 0.$$

The transportation distance $\text{Tran}(w, \alpha; \tilde{w}, \tilde{\alpha})$ between a k -spike distribution $(w, \alpha = (\alpha_1, \dots, \alpha_k))$ and an ℓ -spike distribution $(\tilde{w}, \tilde{\alpha} = (\tilde{\alpha}_1, \dots, \tilde{\alpha}_\ell))$ is defined as the optimum value of the above LP with the objective function replaced by $\sum_{i \in [k], j \in [\ell]} x_{ij} |\alpha_i - \tilde{\alpha}_j|$. Observe that if we view (w, α) equivalently as a k -mixture source $(w, (f^1, \dots, f^k))$ on $\{0, 1\}$ with $f_1^t = \alpha_t$, and $(\tilde{w}, \tilde{\alpha})$ similarly as an ℓ -mixture source on $\{0, 1\}$, then $\text{Tran}(w, \tilde{\alpha}; \tilde{w}, \tilde{\alpha})$ is simply the transportation distance between these k - and ℓ -mixture sources on $\{0, 1\}$.

2.3 Perturbation results and operator norm of random matrices

Definition 2.1. The *operator norm* of A (induced by the ℓ_2 norm) is defined by $\|A\|_{\text{op}} = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}$. The *Frobenius norm* of $A = (A_{i,j})$ is defined by $\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}$.

Lemma 2.2 (Weyl; see Theorem 4.3.1 in [27]). *Let A and B be $n \times n$ matrices such that $\|A - B\|_{\text{op}} \leq \rho$. Let $\lambda_1(A) \geq \dots \geq \lambda_n(A)$, and $\lambda_1(B) \geq \dots \geq \lambda_n(B)$ be the sorted list of eigenvalues of A and B respectively. Then $|\lambda_i(A) - \lambda_i(B)| \leq \rho$ for all $i = 1, \dots, n$.*

Lemma 2.3. *Let A, B be $n \times n$ positive semi-definite (PSD) matrices whose nonzero eigenvalues are at least $\varepsilon > 0$. Let Π_A and Π_B be the projection operators onto the subspaces spanned by the eigenvectors of A and B respectively having nonzero eigenvalues. Let $\|A - B\|_{\text{op}} \leq \rho$. Then $\|\Pi_A - \Pi_B\|_{\text{op}} \leq \sqrt{4\rho/\varepsilon}$.*

Proof. Note that $A\Pi_A = A$, $\Pi_A^2 = \Pi_A$, $B\Pi_B = B$, and $\Pi_B^2 = \Pi_B$. Let x be a unit vector. Then $\|(A - B)x\| \leq \rho$ and, since Π_B is a contraction, $\|(A - B)\Pi_B x\| \leq \rho \|\Pi_B x\| \leq \rho$. Now note that $(A - B)\Pi_B x = A\Pi_B x - Bx$ so by the triangle inequality, we have $\|A\Pi_B x - Ax\| \leq 2\rho$. Now we

can also write $A\Pi_Bx - Ax = A(\Pi_B - \Pi_A)x = A(\Pi_A\Pi_B - \Pi_A)x$. Since A here is acting on a vector that has already been projected down by Π_A , we can conclude

$$2\rho\|A\Pi_Bx - Ax\| = \|A(\Pi_A\Pi_B - \Pi_A)x\| \geq \varepsilon\|(\Pi_A\Pi_B - \Pi_A)x\|.$$

Thus, $2\rho/\varepsilon \geq \|(\Pi_A - \Pi_A\Pi_B)x\|$. By the symmetric argument we also can write $2\rho/\varepsilon \geq \|(\Pi_B - \Pi_B\Pi_A)x\|$. Adding these and applying the triangle inequality we have

$$4\rho/\varepsilon \geq \|(\Pi_A - \Pi_A\Pi_B + \Pi_B - \Pi_B\Pi_A)x\| = \|(\Pi_A^2 - \Pi_A\Pi_B - \Pi_B\Pi_A + \Pi_B^2)x\| = \|(\Pi_A - \Pi_B)^2x\|. \quad \blacksquare$$

Theorem 2.4 ([40]). *For every $\mu > 0$, there is a constant $\kappa = \kappa(\mu) = O(\mu) > 0$ such that the following holds. Let $X_{i,j}, 1 \leq i \leq j \leq n$ be independent random variables with $|X_{ij}| \leq K$, $E[X_{i,j}] = 0$, and $\text{Var}(X_{i,j}) \leq \sigma^2$ for all $i, j \in [n]$, where $\sigma \geq \kappa^2 n^{-1/2} K \ln^2 n$. Let A be the symmetric matrix with entries $A_{i,j} = X_{\min(i,j), \max(i,j)}$ for all $i, j \in [n]$. Then, $\Pr[\|A\|_{\text{op}} \leq 2\sigma\sqrt{n} + \kappa(K\sigma)^{1/2}n^{1/4}\ln n] \geq 1 - n^{-\mu}$.*

3 Our algorithm

We now describe our algorithm that uses 1-, 2-, and $(2k-1)$ -snapshots from the mixture source (w, P) . Given a matrix Z , we use $\text{Span}(Z)$ to denote the column space of Z . Let $r = \sum_{t=1}^k w_t p^t$ denote the 1-snapshot distribution of (w, P) . Let M be the $n \times n$ symmetric matrix representing the 2-snapshot distribution of (w, P) ; so $M_{i,j}$ is the probability of obtaining the 2-snapshot $(i, j) \in [n]^2$. Let $R = rr^\dagger$.

Proposition 3.1. $M = \sum_{t=1}^k w_t p^t p^{t\dagger} = R + A$, where $A = \sum_{t=1}^k w_t (p^t - r)(p^t - r)^\dagger$.

Note that M and A are both PSD. We say that (w, P) is ζ -wide if (i) $\|p - q\|_2 \geq \frac{\zeta}{\sqrt{n}}$ for any two distinct $p, q \in P$; and (ii) the smallest non-zero eigenvalue of A is at least $\frac{\zeta^2}{n}$. (Note that (i) holds if $\min_{p, q \in P, p \neq q} \|p - q\|_1 \geq \zeta$.) We assume that $w_{\min} := \min_t w_t > 0$. Let $k' = \text{rank}(A) \leq k - 1$. It is easy to estimate r using Chernoff bounds (see, e.g., [2]).

Lemma 3.2. *For every $\mu \in \mathbb{N}$ and every $\sigma > 0$, if we use $N \geq \frac{8(\mu+2)}{\sigma^3} \cdot n \ln n$ independent 1-snapshots and set \tilde{r}_i to be the frequency of i in these 1-snapshots for all $i \in [n]$, then with probability at least $1 - n^{-\mu}$ the following hold.*

$$(1 - \sigma)r_i \leq \tilde{r}_i \leq (1 + \sigma)r_i \quad \forall i \text{ with } r_i \geq \frac{\sigma}{2n}, \quad \tilde{r}_i \leq (1 + \sigma)\sigma/2n \quad \forall i \text{ with } r_i < \frac{\sigma}{2n}. \quad (2)$$

It will be convenient in the sequel to assume that our mixture source (w, P) is *isotropic*, by which we mean that $\frac{1}{2n} \leq r_i \leq \frac{2}{n}$ for all $i \in [n]$; notice that this implies that $p_i^t \leq \frac{2}{w_{\min}n}$ for all $i \in [n]$. We show below that this can be assumed at the expense of a small additive error.

Lemma 3.3. *Given an estimate \tilde{r} satisfying (2), if we can learn the constituents of an isotropic k -mixture source to within transportation distance ϵ , then we can learn the constituents of the original k -mixture source to within transportation distance $\epsilon + 4\sigma$.*

Proof. Given an arbitrary mixture source for which we have computed the estimate \tilde{r} , consider the following modification of the mixture constituents. Let $\sigma < 1/32$. We eliminate items i such that $\tilde{r}_i < \frac{2\sigma}{n}$. Each remaining item i is “split” into $\lfloor n\tilde{r}_i/\sigma \rfloor$ items, and the probability of i is split equally among its copies. We can sample m -snapshots from the modified mixture source as follows. We eliminate snapshots that include an eliminated item. With probability at least $1 - n^{-\mu}$, we have

$r_i < \frac{4\sigma}{n}$ if i is eliminated, and $\frac{31}{32} \leq \frac{\tilde{r}_i}{r_i} \leq \frac{33}{32}$ otherwise. So the total weight of eliminated items is at most 4σ and the probability that an m -snapshot survives is at least $(1 - 4\sigma)^m$; we can take $\sigma \ll \frac{1}{m}$. (Recall that we are aiming for $m \leq 2k - 1$.) In the surviving snapshots, each item i in the original snapshot is replaced by one of its $\lfloor n\tilde{r}_i/\sigma \rfloor$ copies, chosen uniformly at random (and independently of previous such choices).

Abusing notation, let p^1, p^2, \dots, p^k be the modified mixture constituents, and r denote the distribution of the modified 1-snapshots. But \tilde{r} still refers to the original mixture.

It is easy to show that the modified mixture source is isotropic. The number n' of new items is at most $\frac{n}{\sigma}$ and at least $\sum_{i: \tilde{r}_i \geq 2\sigma/n} \frac{2}{3} \cdot \frac{n\tilde{r}_i}{\sigma} \geq \frac{2}{3} \cdot \frac{n}{\sigma} \cdot (1 - 2\sigma) \geq \frac{5n}{8\sigma}$. Let $K = \sum_{i: \tilde{r}_i \geq 2\sigma/n} r_i \geq 1 - 4\sigma \geq 7/8$. With probability at least $1 - n^{-\mu}$, for every new item i_ℓ obtained by splitting item i , we have $r_{i_\ell} \geq \frac{r_i}{n\tilde{r}_i/\sigma} \geq \frac{32\sigma}{33n} \geq \frac{1}{2n'}$ and $r_{i_\ell} \leq \frac{1}{K} \cdot \frac{3}{2} \cdot \frac{r_i}{n\tilde{r}_i/\sigma} \leq \frac{384\sigma}{217n} \leq \frac{2}{n'}$. Note that letting $\tilde{r}_{i_\ell} = \tilde{r}_i / \lfloor n\tilde{r}_i/\sigma \rfloor$, we have that $(1 - \sigma)r_{i_\ell} \leq \tilde{r}_{i_\ell} \leq (1 + \sigma)r_{i_\ell}$; thus, we immediately have a good estimate of r for the modified mixture source.

If we learn the modified mixture source, we can revert to the original mixture source by aggregating for each constituent of the mixture the probabilities of the items that we split, and setting the probability of eliminated items to 0. This degrades the quality of the solution by the weight of the eliminated items, which is at most an additive 4σ term in the transportation distance. ■

An overview. Our algorithm for learning an isotropic k -mixture source on $[n]$ takes three parameters: $\zeta \leq 1$ such that (w, P) is ζ -wide, $\omega \in \mathbb{N}$, which controls the success probability of the algorithm, and $\delta \in (0, 1)$, which controls the statistical distance between the constituents of the learnt model and the constituents of the correct model. For convenience, we assume that δ is sufficiently small. The output of the algorithm is a k -mixture source (\tilde{w}, \tilde{P}) such that with probability $1 - O(\frac{1}{\omega})$, $\|w - \tilde{w}\|_\infty$ and $\|p^t - \tilde{p}^t\|_1$ for all $t \in [k]$ tend to 0 as $\delta \rightarrow 0$ (see Theorem 4.1).

The algorithm (see Algorithm 1) consists of three stages. First, we reduce the dimensionality of the problem from n to k' using only 1- and 2-snapshots. By Lemma 3.2, we have an estimate \tilde{r} that is component-wise close to r . Thus, $\tilde{R} = \tilde{r}\tilde{r}^\dagger$ is close in operator norm to R . So we focus on learning the column space of A for which we employ spectral techniques. Leveraging Theorem 2.4, we argue (Lemma 4.2) that by using $O(n \ln^6 n)$ 2-snapshots, one can compute (with high probability) a good enough estimate \tilde{M} of M , and hence obtain a PSD matrix \tilde{A} such that $\|A - \tilde{A}\|_{\text{op}}$ is small.

The remaining task is to learn the projection of P on the affine space $\tilde{r} + \text{Span}(\tilde{A})$, and the mixture weights, which then yields the desired k -mixture source (\tilde{w}, \tilde{P}) . We divide this task into two steps. We choose a random orthonormal basis $\{b_1, \dots, b_{k'}\}$ of $\text{Span}(\tilde{A})$. For each b_j , we consider the projected k -mixture source $(w, \pi_{b_j}(P))$ on \mathbb{R} . One of our contributions is a procedure we devise in Section 5 to learn the corresponding k -spike distribution $(w, \mathbb{E}[\pi_{b_j}(P)])$ using $(2k - 1)$ -snapshots from $(w, \pi_{b_j}(P))$ (which one can easily obtain using $(2k - 1)$ -snapshots from (w, P)). Applying this procedure (see Lemma 4.7), we obtain weights $\tilde{w}_1^j, \dots, \tilde{w}_k^j$ and k (distinct) values $\alpha_1^j, \dots, \alpha_k^j$ such that each true spike $(w_t, b_j^\dagger p^t)$ maps to a distinct inferred spike $(\tilde{w}_{\sigma^j(t)}^j, \alpha_{\sigma^j(t)}^j)$.

The final step is to match up σ_j and $\sigma_{k'}$ for all $j = 1, \dots, k' - 1$, and thus obtain k points in $\tilde{r} + \text{Span}(\tilde{A})$ that are close to the projection of P on $\tilde{r} + \text{Span}(\tilde{A})$. For every $j = 1, \dots, k' - 1$, we generate a random unit “test vector” z_j in $\text{Span}(b_j, b_{k'})$ and learn the projections $z_j^\dagger p^1, \dots, z_j^\dagger p^k$. Since (w, P) is ζ -wide, using standard results about random projections, and the guarantees obtained from our k -spike learning procedure, we can argue that $z_j^\dagger(\alpha_{t_1}^j b_j + \alpha_{t_2}^{k'} b_{k'})$ is close to some value in $\{z_j^\dagger p^t\}_{t \in [k]}$ iff there is some t such that $\alpha_{t_1}^j$ and $\alpha_{t_2}^{k'}$ are close respectively to $b_j^\dagger p_t$ and $b_{k'}^\dagger p_t$ (see Lemma 4.8). Thus, we can use the learned projections of $\{z_j^\dagger p^t\}_{t \in [k]}$ to match up $\{\alpha_t^j\}_{t \in [k]}$ and $\{\alpha_t^{k'}\}_{t \in [k]}$.

Algorithm 1. *Input: an isotropic ζ -wide k -mixture source (w, P) on $[n]$, and parameters $\omega > 0$ and $\delta > 0$. Output: a k -mixture source (\tilde{w}, \tilde{P}) on $[n]$ that is “close” to (w, P) .*

Define $T = 3\omega k^4$, $H = \frac{4}{w_{\min}^2 \zeta \sqrt{n}}$ and $L = \frac{\zeta}{64\omega^{1.5} k^4 \sqrt{n}}$. We assume that $\delta \leq \frac{w_{\min}^3 \zeta^4}{2^{29} \omega^5 k^{16}}$. Let $\kappa = \kappa(2 + \ln \omega)$ be given by Theorem 2.4; we assume $\kappa \geq 1$ for convenience. Define $c = \frac{6400\kappa^2}{w_{\min}^2 \delta^2} \cdot \ln(\frac{1}{\delta})$. We assume that $w_{\min}^2 \geq \frac{240\kappa \ln^{2.5} n}{\sqrt{n}}$.

A1. Dimension reduction.

- A1.1 Use Lemma 3.2 with $\mu = 2 + \ln \omega$ and $\sigma = \frac{\delta}{48}$ to compute an estimate \tilde{r} of r . Set $\tilde{R} = \tilde{r}\tilde{r}^\dagger$.
- A1.2 Independent of all other random variables, choose a Poisson random variable N_2 with expectation $\mathbb{E}[N_2] = cn \ln^6 n$. Choose N_2 independent 2-snapshots and construct a symmetric $n \times n$ matrix \tilde{M} as follows: set $\tilde{M}_{i,i}$ = frequency of the 2-snapshot (i, i) in the sample for all $i \in [n]$, and $\tilde{M}_{i,j} = \tilde{M}_{j,i}$ = half the combined frequency of 2-snapshots (i, j) and (j, i) in the sample, for all $i, j \in [n], i \neq j$.
- A1.3 Compute the spectral decomposition $\tilde{M} - \tilde{R} = \sum_{i=1}^n \lambda_i v_i v_i^\dagger$ where $\lambda_1 \geq \dots \geq \lambda_n$.
- A1.4 Set $\tilde{A} = \sum_{i: \lambda_i \geq \zeta^2/2n} \lambda_i v_i v_i^\dagger$. Note that \tilde{A} is PSD.

A2. Learning projections of (w, P) on random vectors in $\text{Span}(\tilde{A})$.

- A2.1 Pick an orthonormal basis $B = \{b_1, \dots, b_{k'}\}$ for $\text{Span}(\tilde{A})$ uniformly at random.
- A2.2 Set $(\tilde{w}^j, \hat{\alpha}^j) \leftarrow \text{Learn}(b_j, \delta, \frac{1}{6\omega k})$ for all $j = 1, \dots, k'$.

A3. Combining the projections to obtain (\tilde{w}, \tilde{P}) .

- A3.1 Pick $\theta \in [0, 2\pi]$ uniformly at random.
- A3.2 For each $j = 1, \dots, k' - 1$, we do the following.
 - Set $z_j = b_j \cos \theta + b_{k'} \sin \theta$.
 - Set $(\hat{w}^j, \hat{\alpha}^j) \leftarrow \text{Learn}(z_j, \delta, \frac{1}{6\omega k})$.
 - For each $t_1, t_2 \in [k]$, if there exists $t \in [k]$ such that $|(\alpha_{t_1}^j b_j + \alpha_{t_2}^{k'} b_{k'})^\dagger z_j - \hat{\alpha}_t^j| \leq (\sqrt{2} + 1)L/(2 + 5T)$ then set $\varrho^j(t_2) = t_1$.
- A3.3 Define $\varrho^{k'}(t) = t$ for all $t \in [k]$.
- A3.4 For every $t \in [k]$: define $\tilde{w}_t = (\sum_{j=1}^{k'} \tilde{w}_{\varrho^j(t)}^j)/k'$; define $\hat{p}^t = \tilde{r} + \sum_{j=1}^{k'} (\alpha_{\varrho^j(t)}^j - b_j^\dagger \tilde{r}) b_j$, and let \tilde{p}^t be the point in Δ^{n-1} closest in ℓ_1 norm to \hat{p}^t (which can be computed by solving an LP). Return $(\tilde{w}, \tilde{P} = (\tilde{p}^1, \dots, \tilde{p}^k))$.

Algorithm Learn($v, \varsigma, \varepsilon$)

Input: a unit vector $v \in \text{Span}(\tilde{A})$, and parameters $\varsigma > 0, \varepsilon > 0$. We assume that (a) $|v^\dagger(p - q)| \geq L$ for all distinct $p, q \in P$; and (b) $1024k\varsigma < \frac{w_{\min} L}{16H}$.

Output: a k -spike distribution $(\bar{w}, (\gamma_1, \dots, \gamma_k))$ close to $(w, \mathbb{E}[\pi_v(P)])$.

L1. Solve the minimization problem

$$\text{minimize } \|x\|_\infty \quad \text{s.t.} \quad v^\dagger x \geq 1 - \frac{4\delta}{\zeta^2}, \quad \|x\|_2 \leq 1 \quad (\text{Q}_v)$$

which can be formulated as a convex program, to obtain a vector x^* ; set $a = \frac{x^*}{\|x^*\|_2}$. We prove in Lemma 4.4 that $\|a\|_\infty \leq H$ and $|a^\dagger(p - q)| \geq L/2$ for every two mixture constituents $p, q \in P$.

L2. Let $s = \varsigma^{4k}$. Apply the procedure in Section 5 leading to Theorem 5.1 for $(w, \pi_{a/2H}(P))$ to infer a k -spike distribution (\bar{w}, β) that, with probability at least $1 - \varepsilon$, is within transportation distance $O(s^{\Omega(1/k)})$ from $(w, \mathbb{E}[\pi_{a/2H}(P)])$. This uses a sample of $(2k - 1)$ -snapshots of size $3k2^{4k}s^{-4k} \ln(4k/\varepsilon)$.

L3. For every $t \in [k]$, set $\gamma_t = (2H\beta_t)(a^\dagger v)$. Return (\bar{w}, γ) .

4 Analysis

We prove the following theorem.

Theorem 4.1. *Algorithm 1 uses $O(\frac{\ln \omega}{\delta^3} \cdot n \ln n)$ 1-snapshots, $O(\frac{\ln^2 \omega \ln(1/\delta)}{\delta^2 w_{\min}^2} \cdot n \ln^6 n)$ 2-snapshots, and $O(\frac{k 2^{4k}}{\delta^{16k^2}} \cdot \ln(24\omega k^2))$ $(2k-1)$ -snapshots, and computes a k -mixture source (\tilde{w}, \tilde{P}) on $[n]$ such that with probability $1 - O(\frac{1}{\omega})$, there is a permutation $\sigma : [k] \mapsto [k]$ such that for all $t = 1, \dots, k$,*

$$|w_t - \tilde{w}_{\sigma(t)}| = O\left(\frac{\delta \omega^{1.5} k^5}{w_{\min}^2 \zeta^2}\right) \quad \text{and} \quad \|p^t - \tilde{p}^{\sigma(t)}\|_1 = O\left(\frac{\sqrt{k\delta}}{w_{\min}^{1.5} \zeta}\right).$$

$$\text{Hence, } \text{Tran}(w, P; \tilde{w}, \tilde{P}) = O\left(\frac{\sqrt{k\delta}}{w_{\min}^{1.5} \zeta}\right).$$

The roadmap of the proof is as follows. By Lemma 3.2, with probability at least $1 - \frac{1}{\omega n^2}$, $(1 - \frac{\delta}{48})r_i \leq \tilde{r}_i \leq (1 + \frac{\delta}{48})r_i$ for all $i \in [n]$. We assume that this holds in the sequel. In Lemma 4.2, we prove that the matrix \tilde{A} computed after step A1 is a good estimate of A . In Lemma 4.3, we derive some properties of the column space of A . Lemma 4.4 then uses these properties to show that algorithm Learn returns a good approximation to $(w, \mathbb{E}[\pi_v(P)])$. Claim 4.5 and Corollary 4.6 prove that the projections of the mixture constituents on the b_j s and the z_j s are well-separated. Combining this with Lemma 4.4, we prove in Lemma 4.7 that with suitably large probability, every true spike $(w_t, b_j^\dagger p^t)$ maps to a distinct nearby inferred spike on every b_j , $j \in [k']$, and similarly every true spike $(w_t, z_j^\dagger p^t)$ maps to a distinct nearby inferred spike on every z_j , $j \in [k' - 1]$. Lemma 4.8 shows that one can then match up the spikes on the different b_j s. This yields k points in $\text{Span}(\tilde{A})$ that are close to the projection of P on $\text{Span}(\tilde{A})$. Finally, we argue that this can be mapped to a k -mixture source (\tilde{w}, \tilde{P}) that is close to (w, P) .

Lemma 4.2. *With probability at least $1 - \frac{1}{n\omega}$, the matrix \tilde{A} computed after step A1 satisfies $\text{rank}(\tilde{A}) = k' = \text{rank}(A)$ and $\|A - \tilde{A}\|_{\text{op}} \leq \frac{\delta}{n}$.*

Proof. Recall that $k' = \text{rank}(A)$. Let $B = \tilde{M} - \tilde{R} = \sum_{i=1}^n \lambda_i v_i v_i^\dagger$, where $\lambda_1 \geq \dots \geq \lambda_n$. We prove below that with probability at least $1 - \frac{1}{n\omega}$, we have $\|M - \tilde{M}\|_{\text{op}} \leq \frac{\delta}{4n}$ and $\|R - \tilde{R}\|_{\text{op}} \leq \frac{\delta}{4n}$. This implies that $\|A - B\|_{\text{op}} \leq \|M - \tilde{M}\|_{\text{op}} + \|R - \tilde{R}\|_{\text{op}} \leq \frac{\delta}{2n}$. Hence, by Lemma 2.2, it follows that by the ζ -wide assumption, $\lambda_{k'} \geq \frac{\zeta^2}{n} - \frac{\delta}{2n} \geq \frac{3\zeta^2}{4n}$, and $|\lambda_i| \leq \frac{\delta}{2n} \leq \frac{\zeta^2}{4n}$ for all $i > k'$. Thus, we include exactly k' eigenvectors when defining \tilde{A} , so $\text{rank}(\tilde{A}) = k'$. Since \tilde{A} is the closest rank- k' approximation in operator norm to B , we have $\|A - \tilde{A}\|_{\text{op}} \leq \|A - B\|_{\text{op}} + \|B - \tilde{A}\|_{\text{op}} \leq 2\|A - B\|_{\text{op}} \leq \frac{\delta}{n}$.

We now proceed to bound $\|M - \tilde{M}\|_{\text{op}}$ and $\|R - \tilde{R}\|_{\text{op}}$. It is easy to see that $|\tilde{R}_{i,j} - R_{i,j}| \leq 3\sigma r_{i,j}$, where $\sigma = \delta/48$, and so $\|R - \tilde{R}\|_{\text{op}} \leq \|R - \tilde{R}\|_F \leq \frac{\delta}{4n}$.

Bounding $\|M - \tilde{M}\|_{\text{op}}$ is more challenging. Note that $M_{i,j} \leq \min\{\frac{2}{n}, \frac{4}{w_{\min} n^2}\}$ due to isotropy. Let $K = \frac{4 \ln(1/\delta)}{\delta}$. Let $D = N_2 \cdot (\tilde{M} - M)$. Let $X_{i,i}^\ell = 1$ if the ℓ -th snapshot is (i, i) , for $i \in [n]$, and for $i, j \in [n], i \neq j$, let $X_{i,j}^\ell = X_{j,i}^\ell = \frac{1}{2}$ if the ℓ -th 2-snapshot is (i, j) or (j, i) , and 0 otherwise. Let $Y_{i,j}^\ell = X_{i,j}^\ell - M_{i,j} = X_{i,j}^\ell - \mathbb{E}[X_{i,j}^\ell]$; so $D_{i,j} = \sum_{\ell=1}^{N_2} Y_{i,j}^\ell$ for all $i, j \in [n]$. We have $\sigma^2(n_2) := \text{Var}[D_{i,j} | N_2 = n_2] = n_2 \text{Var}[X_{i,j}^1] \leq n_2 \mathbb{E}[(X_{i,j}^1)^2] \leq n_2 M_{i,j}$. For $n_2 \leq 2cn \ln^6 n$, we have

$\sigma^2(n_2) \leq \frac{8c \ln^6 n}{w_{\min}^2 n} \leq \frac{\ln n \ln(1/\delta)}{\delta^2}$ (since $w_{\min}^4 \geq \frac{57600 \kappa^2 \ln^5 n}{n}$). So by Bernstein's inequality,

$$\begin{aligned} \Pr[|D_{i,j}| > K \ln n | N_2 = n_2] &\leq 2 \exp\left(-\frac{K^2 \ln^2 n}{2(\sigma^2(n_2) + K \ln n/3)}\right) \\ &\leq 2 \max\left\{\exp\left(-\frac{K^2 \ln^2 n}{4\sigma^2(n_2)}\right), \exp\left(-\frac{3K \ln n}{4}\right)\right\} \leq \frac{2\delta}{n^3}. \end{aligned}$$

Since $\Pr[N_2 > 2c \ln^6 n] \leq n^{-3}$, we can say that with probability at least $1 - 2n^{-2}$, we have $|D_{i,j}| \leq K \ln n$ for every $i, j \in [n]$ and $N_2 \leq 2c \ln^6 n$.

Define a matrix D' by putting, for every $i, j \in [n]$, $D'_{i,j} = \text{sign}(D_{i,j}) \cdot \min\{|D_{i,j}|, K \ln n\}$. Put $D'' = D' - \mathbb{E}[D']$. Clearly, $\mathbb{E}[D''_{i,j}] = 0$ for every $i, j \in [n]$. The entries of D are independent random variables as N_2 is a Poisson random variable; hence, the entries of D'' are also independent random variables. Also $\text{Var}[D''_{i,j}] \leq \text{Var}[D_{i,j}]$ since censoring a random variable to an interval can only reduce the variance. Note that $D_{i,j} = \sum_{\ell=1}^{N_2} Y_{i,j}^\ell$ follows the compound Poisson distribution. So we have

$$\text{Var}[D_{i,j}] = \mathbb{E}[N_2] \cdot \mathbb{E}[(Y_{i,j}^1)^2] = \mathbb{E}[N_2] \cdot \text{Var}[X_{i,j}^1] \leq \mathbb{E}[N_2] M_{i,j} \leq \frac{4c \ln^6 n}{w_{\min}^2 \cdot n} \leq \frac{\hat{c}^2 K^2 \ln^6 n}{n}$$

where $\hat{c} = \max\left\{\frac{2\sqrt{c}}{w_{\min} K}, \kappa^2\right\}$. Thus, by Theorem 2.4, the constant $\kappa = \kappa(2 + \ln \omega) > 0$ is such that with probability at least $1 - \frac{1}{n^2 \omega}$

$$\|D''\|_{\text{op}} \leq 2 \cdot \frac{\hat{c} K \ln^3 n}{\sqrt{n}} \cdot \sqrt{n} + \kappa \sqrt{K \ln n \cdot \frac{\hat{c} K \ln^3 n}{\sqrt{n}}} \cdot \sqrt[4]{n} \cdot \ln n \leq (2K\hat{c} + \kappa K \sqrt{\hat{c}}) \ln^3 n.$$

We have $\Pr[N_2 \geq \frac{1}{2} \mathbb{E}[N_2]] \geq 1 - n^{-2}$. Thus, with probability at least $1 - \frac{1}{n\omega}$, we have that $N_2 \geq \frac{1}{2} \mathbb{E}[N_2]$, $D' = D$, and $\|D''\|_{\text{op}}$ is bounded by (3). We show below that $2\|E[D']\|_{\text{op}}/\mathbb{E}[N_2] \leq 6\delta n^{-2} \leq \delta/20n$. One can verify that $4K\hat{c}/c \leq \delta/10$ and $2\kappa K \sqrt{\hat{c}}/c \leq \delta/10$. Therefore, with probability at least $1 - \frac{1}{n\omega}$, we have that $\|M - \tilde{M}\|_{\text{op}} = \frac{1}{N_2} \cdot \|D\|_{\text{op}} \leq \frac{2}{\mathbb{E}[N_2]} \cdot (\|D''\|_{\text{op}} + \|E[D']\|_{\text{op}}) \leq \frac{\delta}{4n}$.

Finally, we bound $\|E[D']\|_{\text{op}}$. We have $\|E[D']\|_{\text{op}} \leq \|E[D']\|_F = \|E[D' - D]\|_F \leq n \cdot \max_{i,j} \mathbb{E}[|D'_{i,j} - D_{i,j}|]$. Let $\mu = cn \ln^6 n = \mathbb{E}[N_2]$. Fix any i, j . For any $n_2 \leq 2 \ln(1/\delta)\mu$, we have $\text{Var}[D_{i,j} | N_2 = n_2] \leq n_2 M_{i,j} \leq \frac{8c \ln(1/\delta) \ln^6 n}{w_{\min}^2 n}$. So by Bernstein's inequality, we have that $\Pr[|D_{i,j}| > K \ln n | N_2 \leq 2 \ln(1/\delta)\mu] < 2\delta n^{-3}$. Also, $|D'_{i,j} - D_{i,j}| \leq N_2$ always. Therefore,

$$\mathbb{E}[|D'_{i,j} - D_{i,j}| | N_2 = n_2] \leq \begin{cases} 2\delta n^{-3} n_2 & \text{if } n_2 \leq 2 \ln(1/\delta)\mu; \\ n_2 & \text{otherwise} \end{cases}$$

and $\mathbb{E}[|D'_{i,j} - D_{i,j}|] \leq \mu - \Pr[N_2 \leq 2 \ln(1/\delta)\mu] \mathbb{E}[N_2 | N_2 \leq 2 \ln(1/\delta)\mu] (1 - 2\delta n^{-3})$. Since N_2 is Poisson distributed, we have

$$\begin{aligned} \Pr[N_2 \leq 2 \ln(1/\delta)\mu] \mathbb{E}[N_2 | N_2 \leq 2 \ln(1/\delta)\mu] &= \sum_{\ell=0}^{[2 \ln(1/\delta)\mu]} \ell \cdot \frac{\mu^\ell e^{-\mu}}{\ell!} = \mu \sum_{\ell=0}^{[2 \ln(1/\delta)\mu]-1} \frac{\mu^\ell e^{-\mu}}{\ell!} \\ &\geq \mu \Pr[N_2 \leq \ln(1/\delta)\mu] \geq \mu(1 - \delta n^{-3}). \end{aligned}$$

Thus, $\mathbb{E}[|D'_{i,j} - D_{i,j}|] \leq \mu - \mu(1 - \delta n^{-3})(1 - 2\delta n^{-3}) \leq 3\delta n^{-3}\mu$, and $2\|E[D']\|_{\text{op}}/\mathbb{E}[N_2] \leq 6\delta n^{-2}$. ■

We assume in the sequel that the high-probability event stated in Lemma 4.2 happens. Thus, Lemma 2.3 implies that $\|\Pi_A - \Pi_{\tilde{A}}\|_{\text{op}} \leq \frac{2\sqrt{\delta}}{\zeta}$.

Lemma 4.3. For every unit vector $b \in \text{Span}(A)$, $\|b\|_\infty \leq \frac{2}{w_{\min}^2 \zeta \sqrt{n}}$.

Proof. Recall that $A = \sum_{t=1}^k w_t (p^t - r)(p^t - r)^\dagger$, and the smallest non-zero eigenvalue of A is at least ζ^2/n . Note that $\text{Span}(A) = \text{Span}\{p^1 - r, \dots, p^k - r\}$. Let $Z = \text{conv}(P)$. As all the mixture constituents $p \in P$ satisfy $\|p\|_\infty \leq \frac{2}{w_{\min} n}$, for any $x \in Z$, clearly $\|x\|_\infty \leq \frac{2}{w_{\min} n}$. Also since $x \geq 0$ and $\|r\|_\infty \leq \frac{2}{n}$, we have $\|x - r\|_\infty \leq \frac{2}{w_{\min} n}$. So if $r + b \in Z$, then $\|b\|_\infty \leq \frac{2}{w_{\min} n}$. Otherwise, let the line segment $[r, r + b]$ intersect the boundary of Z at some point b' . We show that $\|r - b'\|_2^2 \geq \frac{\zeta^2 w_{\min}^2}{n}$. The lemma then follows since $b = (b' - r)/\|b' - r\|_2$ and so $\|b\|_\infty = \frac{\|b' - r\|_\infty}{\|b' - r\|_2} \leq \frac{2}{w_{\min}^2 \zeta \sqrt{n}}$.

Let S be a facet of Z such that $b' \in S$, $r \notin S$ (there must be such a facet since r is in the strict interior of P as $w_{\min} > 0$). Since $Z \subseteq \text{Span}(A)$, one can find a unit vector $v \in \text{Span}(A)$ such that S is exactly the set of points that minimize $v^\dagger x$ over $x \in Z$. Let $d_L = v^\dagger r - \min_{x \in Z} v^\dagger x = v^\dagger (r - b')$. We lower bound $\|r - b'\|_2$ by d_L . Note that $d_L > 0$. Clearly, $v^\dagger (p^t - r) \geq -d_L$ for all $t \in [k]$. Projecting P onto v , we have that (a) $\sum_{t=1}^k w_t v^\dagger (p^t - r) = 0$; and (b) $v^T A v = \sum_{t=1}^k w_t (v^\dagger (p^t - r))^2 \geq \frac{\zeta^2}{n}$ since $v \in \text{Span}(A)$ and (w, P) is ζ -wide. Let $W_L = \sum_{t: v^\dagger (p^t - r) \leq 0} w_t$, let $W_R = 1 - W_L \geq w_{\min}$, and let $d_R = \max_t \{v^\dagger (p^t - r)\}$. Then, $0 = \sum_{t=1}^k w_t v^\dagger (p^t - r) \geq W_L(-d_L) + w_{\min} d_R$, so $d_R \leq d_L \cdot \frac{W_L}{w_{\min}}$. Also $\frac{\zeta^2}{n} \leq \sum_{t=1}^k w_t (v^\dagger (p^t - r))^2 \leq W_L \cdot d_L^2 + W_R \cdot d_R^2 \leq W_L \cdot d_L^2 + W_R \cdot d_L^2 \cdot \frac{W_L^2}{w_{\min}^2} \leq \frac{d_L^2}{w_{\min}^2}$. So, $d_L^2 \geq \frac{\zeta^2 w_{\min}^2}{n}$. \blacksquare

Lemma 4.4. If the assumptions stated in Algorithm Learn are satisfied, then the vector a computed in Learn satisfies $\|a\|_\infty \leq H$, and $|a^\dagger(p - q)| \geq L/2$ for every two mixture constituents $p, q \in P$. Hence, with probability at least $1 - \varepsilon$, the output (\bar{w}, γ) of Learn satisfies the following: there is a permutation $\sigma : [k] \mapsto [k]$ such that for all $t = 1, \dots, k$,

$$|w_t - \bar{w}_{\sigma(t)}| = O\left(\frac{\zeta \omega^{1.5} k^5}{w_{\min}^2 \zeta^2}\right), \quad |v^\dagger p^t - \gamma_{\sigma(t)}| \leq \frac{2048kH\zeta}{w_{\min}} + \frac{8\sqrt{2\delta}}{w_{\min}\zeta\sqrt{n}} \leq \frac{2048kH\zeta}{w_{\min}} + \frac{L}{8T}.$$

Proof. We have $v^\dagger \Pi_A(v) = 1 - \|v - \Pi_A(v)\|_2^2 = 1 - \|(\Pi_{\bar{A}} - \Pi_A)v\|_2^2 \geq 1 - \frac{4\delta}{\zeta^2}$. Thus, $\Pi_A(v)$ is feasible to (Q_v) , and since $\|\Pi_A(v)\|_2 \leq 1$, by Lemma 4.3, the optimal solution x^* to (Q_v) satisfies $\|x^*\|_\infty \leq \|\Pi_A(v)\|_\infty \leq H/2$. Also $\|x^*\|_2^2 \geq v^\dagger x^* \geq 1 - \frac{4\delta}{\zeta^2} \geq \frac{1}{4}$, so $\|a\|_\infty \leq H$. Note that $\|v - a\|_2^2 = 2(1 - v^\dagger a) \leq 2(1 - v^\dagger x^*) \leq \frac{8\delta}{\zeta^2}$. It follows that for any two mixture constituents p, q , we have

$$\begin{aligned} |a^\dagger(p - q)| &\geq |v^\dagger(p - q)| - |(v - a)^\dagger(p - q)| \geq |v^\dagger(p - q)| - \frac{2\sqrt{2\delta}}{\zeta} \|p - q\|_2 \\ &\geq |v^\dagger(p - q)| - \frac{8\sqrt{2\delta}}{w_{\min}\zeta\sqrt{n}} \geq |v^\dagger(p - q)| - \frac{L}{2} \geq \frac{L}{2}. \end{aligned}$$

So any two spikes in the k -spike mixture $(w, \mathbb{E}[\pi_{a/2H}(P)])$ are separated by a distance of at least $L/4H$. Since $s < L/4H$, Theorem 5.1 guarantees that with a sample of $(2k - 1)$ -snapshots of size $3k2^{4k}s^{-4k} \log(2k/\varepsilon)$, with probability at least $1 - \varepsilon$, the learned k -spike distribution (\bar{w}, β) satisfies $\text{Tran}(w, \mathbb{E}[\pi_{a/2H}(P)]; \bar{w}, \beta) \leq 1024ks^{1/(4k)} = 1024k\zeta < \frac{Lw_{\min}}{8H}$. Notice that this implies that there is a permutation $\sigma : [k] \mapsto [k]$ such that $\forall t = 1, \dots, k$:

$$|(a/2H)^\dagger p^t - \beta_{\sigma(t)}| \leq \frac{1024k\zeta}{w_{\min}} < \frac{L}{8H}, \quad |w_t - \bar{w}_{\sigma(t)}| = O\left(\frac{k\zeta}{L/8H}\right) = O\left(\frac{\zeta \omega^{1.5} k^5}{w_{\min}^2 \zeta^2}\right). \quad (3)$$

Fix some $t \in [k]$. Let $t' = \sigma(t)$. From (3), we know that $|a^\dagger p^t - 2H \cdot \beta_{t'}| = \frac{2048kH\zeta}{w_{\min}}$. We bound $|v^\dagger p^t - a^\dagger p^t|$ and $|2H\beta_{t'} - \gamma_{t'}|$, which together with the above will complete the proof of the lemma.

We have $|(v - a)^\dagger p^t| \leq \|v - a\|_2 \|p^t\|_2 \leq \frac{4\sqrt{2\delta}}{w_{\min}\zeta\sqrt{n}}$. Since $\gamma_{t'} = (2H\beta_{t'})a^\dagger v$ and $|\beta_{t'}| \leq \frac{1}{2}$, we have $|2H\beta_{t'} - \gamma_{t'}| \leq \frac{H\cdot 4\delta}{\zeta^2} \leq \frac{16\delta}{w_{\min}^2\zeta^3\sqrt{n}}$. It follows that $|v^\dagger p^t - \gamma_{t'}| \leq \frac{2048kH\zeta}{w_{\min}} + \frac{8\sqrt{2\delta}}{w_{\min}\zeta\sqrt{n}} \leq \frac{2048kH\zeta}{w_{\min}} + \frac{L}{8T}$. ■

Claim 4.5. *Let Z be a random unit vector in $\text{Span}(\tilde{A})$ and $v \in \text{Span}(\tilde{A})$. $\Pr[|Z^\dagger v| < \frac{\|v\|_2}{20\omega^{1.5}k^4}] < \frac{1}{3\omega k'k^2}$.*

Proof. One way of choosing the random unit vector Z is as follows. Fix an orthonormal basis $\{u_1, \dots, u_{k'}\}$ for $\text{Span}(\tilde{A})$. We choose independent $N(0, 1)$ random variables X_i for $i \in [k']$. Define $C = \sum_{i=1}^{k'} X_i u_i$ and set $Z = C/\|C\|_2$. Set $a_1 = \frac{\pi}{32\omega^2 k'^2 k^4}$ and $a_2 = 2 + \frac{4\ln(12\omega k'k^2)}{k'} \leq 96\omega k$.

Note that $C^\dagger v/\|v\|_2$ is distributed as $N(0, 1)$. Therefore, $\Pr[|C^\dagger v| \leq \|v\|_2 \sqrt{a_1}] \leq \sqrt{\frac{2a_1}{\pi}} \leq \frac{1}{4\omega k'k^2}$. Also, $\|C\|_2^2 = \sum_{i=1}^{k'} X_i^2$ follows the $\chi_{k'}^2$ distribution. So

$$\Pr[\|C\|_2^2 > a_2 k'] < \left(a_2 e^{1-a_2}\right)^{k'/2} < \exp\left((1 - a_2/2)k'/2\right) < \frac{1}{12\omega k'k^2}.$$

Observe that $\sqrt{\frac{a_1}{a_2 k'}} \geq \frac{1}{32\omega^{1.5}k^4}$. So if the “bad” event stated in the lemma happens, then $|C^\dagger v| \leq \|v\|_2 \sqrt{a_1}$ or $\|C\|_2^2 \geq a_2 k'$ happens; the probability of this is at most $\frac{1}{3\omega k'k^2}$. ■

Lemma 4.6. *With probability at least $1 - \frac{1}{3\omega}$, for every pair $p, q \in P$, we have (i) $|b_j^\dagger(p - q)| \geq L$ for every $j \in [k']$ and (ii) $|z_j^\dagger(p - q)| \geq L$ for every $j \in [k' - 1]$.*

Proof. Define $\tilde{p} = \Pi_{\tilde{A}}(p)$ for a mixture constituent p . Clearly, for any $v \in \text{Span}(\tilde{A})$, $v^\dagger \tilde{p} = v^\dagger p$. Recall that $\|\Pi_A - \Pi_{\tilde{A}}\| \leq \frac{2\sqrt{\delta}}{\zeta}$. So for every $p, q \in P$, $\|\tilde{p} - \tilde{q}\|_2^2 \geq \|p - q\|_2^2 - \|(\Pi_A - \Pi_{\tilde{A}})(p - q)\|_2^2 \geq \|p - q\|_2^2/4$; hence, $\|\tilde{p} - \tilde{q}\|_2 \geq \frac{\zeta}{2\sqrt{n}}$. Notice that the z_j vectors are also random unit vectors in $\text{Span}(\tilde{A})$. Applying Claim 4.5 to each event involving one of the $\{b_j\}_{j \in [k']}$, $\{z_j\}_{j \in [k'-1]}$ random unit vectors, and one of the $\binom{k}{2}$ vectors $\|\tilde{p} - \tilde{q}\|$ for $\tilde{p}, \tilde{q} \in \Pi_{\tilde{A}}(P)$, and taking the union bound over the at most $k'k^2$ such events completes the proof. ■

Lemma 4.7. *With probability at least $1 - \frac{2}{3\omega}$, the k -spike distributions obtained in steps A2 and A3 satisfy:*

(i) *For every $j \in [k']$, there is a permutation $\sigma^j : [k] \mapsto [k]$ such that for all $t \in [k]$,*

$$|w_t - \tilde{w}_{\sigma^j(t)}^j| = O\left(\frac{\delta\omega^{1.5}k^5}{w_{\min}^2\zeta^2}\right), \quad |b_j^\dagger p^t - \alpha_{\sigma^j(t)}^j| = O\left(\frac{\sqrt{\delta}}{w_{\min}^{1.5}\zeta\sqrt{n}}\right) \text{ and is at most } \frac{L}{2 + 5T}.$$

Hence, $|\alpha_{t_1}^j - \alpha_{t_2}^j| \geq \frac{L}{1+0.4/T}$ for all distinct $t_1, t_2 \in [k]$.

(ii) *For every $j \in [k' - 1]$, for every $t \in [k]$, there is a distinct t' such that*

$$|w_t - \hat{w}_{t'}^j| = O\left(\frac{\delta\omega^{1.5}k^5}{w_{\min}^2\zeta^2}\right), \quad |z_j^\dagger p^t - \hat{\alpha}_{t'}^j| = O\left(\frac{\sqrt{\delta}}{w_{\min}^{1.5}\zeta\sqrt{n}}\right) \text{ and is at most } \frac{L}{2 + 5T}.$$

Proof. Assume that the event stated in Lemma 4.6 happens. Then the inputs to Learn in steps A2 and A3 are “valid”, i.e., satisfy the assumptions stated in Algorithm Learn. Plug in $\varsigma = \delta$ and $\varepsilon = \frac{1}{6\omega k}$ in Lemma 4.4. Taking the union bound over all the b_j s and the z_j s, we obtain that the probability that Learn fails on some input, when all the b_j s and z_j s are valid is at most $\frac{1}{3\omega}$. The lemma follows from Lemma 4.4 by noting that $\frac{2048kH\delta}{w_{\min}} = O\left(\frac{\sqrt{\delta}}{w_{\min}^{1.5}\zeta\sqrt{n}}\right)$ and is at most $\frac{L}{24T}$, and $L/24T + L/8T \leq L/(2 + 5T)$. ■

Lemma 4.8. *With probability at least $1 - \frac{1}{\omega}$, for every $j = 1, \dots, k' - 1$ ϱ^j is a well-defined function and $\varrho^j(\sigma^j(t)) = \sigma^{k'}(t)$ for every $t \in [k]$.*

Proof. Assume that the events in Lemmas 4.6 and 4.7 occur. Fix $j \in [k' - 1]$. We call a point $\alpha_{t_1}^j b_j + \alpha_{t_2}^{k'} b_{k'}$ a grid- j point. Call this grid point “genuine” if there exists $t \in [k]$ such that $\sigma^j(t) = t_1$ and $\sigma^{k'}(t) = t_2$, and “fake” otherwise. The distance between any two grid- j points is at least $L/(1 + 0.4/T)$ (by Lemma 4.7). So the probability there is a pair of genuine and fake grid- j points whose projections on z_j are less than $L/(T + 0.4)$ away is at most $k^3 \cdot \frac{2}{\pi} \arcsin(\frac{1}{T}) \leq k^3 \cdot \frac{2}{\pi} \cdot \frac{6}{5T} \leq \frac{1}{3\omega k}$. Therefore, with probability at least $1 - \omega$, the events in Corollary 4.6 and Lemma 4.7 happen, and for all $j \in [k' - 1]$, every pair of genuine and fake grid- j points project to points on z_j that are at least $L/(T + 0.4)$ apart. We condition on this in the sequel.

Now fix $j \in [k' - 1]$ and consider any pair $t_1, t_2 \in [k]^2$. Let g be the grid- j point $b_j \alpha_{t_1}^j + b_{k'} \alpha_{t_2}^{k'}$. We show that $\varrho^j(t_2) = t_1$ iff g is a genuine grid- j point. If g is genuine, let t be such that $\sigma^j(t) = t_1$, $\sigma^{k'}(t) = t_2$. Let p^t be the projection of p^t on $\text{Span}(b_j, b_{k'})$. By Lemma 4.7, we have that $\|p' - g\|_2 \leq \frac{\sqrt{2}L}{2+5T}$. Also, there exists $t' \in [k]$ such that $|\hat{\alpha}_{t'}^j - z_j^\dagger p^t| \leq \frac{L}{2+5T}$. Since $z_j^\dagger p' = z_j^\dagger p^t$, this implies that $|z_j^\dagger g - \hat{\alpha}_{t'}^j| \leq |\hat{\alpha}_{t'}^j - z_j^\dagger p^t| + |z_j^\dagger (p' - g)| \leq \frac{(\sqrt{2}+1)L}{2+5T}$ and so $\varrho^j(t_2) = t_1$.

Now suppose g is fake but $|z_j^\dagger g - \hat{\alpha}_{t'}^j| \leq (\sqrt{2} + 1)L/(2 + 5T)$ for some $t' \in [k]$. Let $t \in [k]$ be such that $|\hat{\alpha}_{t'}^j - z_j^\dagger p^t| \leq \frac{L}{2+5T}$. Let g' be the genuine grid point $b_j \alpha_{\sigma^j(t)}^j + b_{k'} \alpha_{\sigma^{k'}(t)}^{k'}$. So $|z_j^\dagger g' - \hat{\alpha}_{t'}^j| \leq (\sqrt{2} + 1)L/(2 + 5T)$, and hence $|z_j^\dagger (g - g')| \leq \frac{2(\sqrt{2}+1)L}{2+5T} < \frac{L}{0.4+T}$ which is a contradiction. \blacksquare

Proof of Theorem 4.1. We condition on the fact that all the “good” events stated in Lemmas 3.2, 4.2, 4.6, 4.7, and 4.8 happen. The probability of success is thus $1 - O(\frac{1}{\omega})$. The sample-size bounds follow from the description of the algorithm. For notational simplicity, let $\sigma^{k'}$ be the identity permutation, i.e., $\sigma^{k'}(t) = t$ for all $t \in [k]$. So by Lemma 4.8, we have $\varrho^j(t) = \sigma^j(t)$ for every $j \in [k' - 1]$ and $t \in [k]$.

For $t = 1, 2, \dots, k$, define $\bar{p}^t = \tilde{r} + \sum_{j=1}^{k'} b_j^\dagger (p^t - \tilde{r}) b_j = \tilde{r} + \Pi_{\tilde{A}}(p^t - \tilde{r})$. Fix $t \in [k]$. Then

$$\|p^t - \tilde{p}^t\|_1 \leq \|p^t - \hat{p}^t\|_1 + \|\hat{p}^t - \tilde{p}^t\|_1 \leq 2\|p^t - \hat{p}^t\|_1 \leq 2(\|p^t - \bar{p}^t\|_1 + \|\bar{p}^t - \tilde{p}^t\|_1).$$

$$\begin{aligned} \text{We have } \|p^t - \bar{p}^t\|_2 &= \left\| r - \tilde{r} + (p^t - r) - \sum_{j=1}^{k'} b_j^\dagger (p^t - \tilde{r}) b_j \right\|_2 \\ &= \left\| r - \tilde{r} + \Pi_{\tilde{A}}(\tilde{r} - r) + (p^t - r) - \Pi_{\tilde{A}}(p^t - r) \right\|_2 \\ &\leq 2 \cdot \|r - \tilde{r}\|_2 + \|\Pi_A - \Pi_{\tilde{A}}\|_{\text{op}} \cdot \|p^t - r\|_2 \leq \frac{\delta}{12\sqrt{n}} + \frac{8\sqrt{2}\delta}{w_{\min}\zeta\sqrt{n}}. \end{aligned}$$

$$\text{Also } \|\bar{p}^t - \hat{p}^t\|_2 \leq \left\| \sum_{j=1}^{k'} (b_j^\dagger p^t - \alpha_{\sigma^j(t)}^j) b_j \right\|_2 = O\left(\frac{\sqrt{k}\delta}{w_{\min}^{1.5}\zeta\sqrt{n}}\right)$$

where the last equality follows from Lemma 4.7. Thus, $\|p^t - \tilde{p}^t\|_1 = O(\frac{\sqrt{k}\delta}{w_{\min}^{1.5}\zeta})$. Also, we have

$|w_t - \tilde{w}_t| = O\left(\frac{\delta \omega^{1.5} k^5}{w_{\min}^2 \zeta^2}\right)$ by Lemma 4.7. Finally, note that

$$\begin{aligned} \text{Tran}(w, P; \tilde{w}, \tilde{P}) &\leq \frac{1}{2} \left(\sum_{t=1}^k \min\{w_t, \tilde{w}_t\} \max_t \|p^t - \tilde{p}^t\|_1 + \|w - \tilde{w}\|_1 \max_{t,t'} \|p^t - \tilde{p}^{t'}\|_1 \right) \\ &\leq \frac{1}{2} \left(\sum_{t=1}^k \min\{w_t, \tilde{w}_t\} \max_t \|p^t - \tilde{p}^t\|_1 + \|w - \tilde{w}\|_1 \left(\max_{t,t'} \|p^t - \tilde{p}^{t'}\|_1 + \max_t \|p^t - \tilde{p}^t\|_1 \right) \right) \\ &\leq \max_t \|p^t - \tilde{p}^t\|_1 + \|w - \tilde{w}\|_1 \cdot \frac{2}{w_{\min}} = O\left(\frac{\sqrt{k\delta}}{w_{\min}^{1.5} \zeta}\right). \quad \blacksquare \end{aligned}$$

5 The one-dimensional problem: learning mixture sources on $[0,1]$

In this section, we supply the key subroutine called upon in step L2 of Algorithm **Learn**, which will complete the description of Algorithm 1. We are given a k -mixture source $(w, \pi_x(P))$ on $[-\frac{1}{2}, \frac{1}{2}]$. (Recall that **Learn** invokes the procedure for the mixture $(w, \pi_{a/2H}(P))$ where $\|a\|_\infty \leq H$.) It is clear that we *cannot* in general reconstruct this mixture source with an aperture size that is independent of n , let alone aperture $2k - 1$. However, our goal is somewhat different and more modest. We seek to reconstruct the k -spike distribution $(w, \mathbb{E}[\pi_x(P)])$, and we show that this *can* be achieved with aperture $2k - 1$ (which is the smallest aperture at which this is information-theoretically possible).

It is easy to obtain a $(2k - 1)$ -snapshot from $(w, \pi_x(P))$ given a $(2k - 1)$ -snapshot from (w, P) by simply replacing each item $i \in [n]$ that appears in the snapshot by x_i . We will assume in the sequel that every constituent $\pi_x(p^t)$ is supported on $[0, 1]$, which is simply a translation by $\frac{1}{2}$.

To simplify notation, we use $\theta = (\vartheta, (q^1, \dots, q^k))$ to denote the k -mixture source on $[0, 1]$, and $(\vartheta, \alpha = (\alpha_1, \dots, \alpha_k))$ to denote the corresponding k -spike distribution, where $\alpha_i \in [0, 1]$ is the expectation of q^i for all $i \in [k]$. We equivalently view (ϑ, α) as a k -mixture source $(\vartheta, (f^1, \dots, f^k))$ on $\{0, 1\}$: each f^i is a “coin” whose bias is $f_1^i = \alpha_i$. In Section 5.1, we describe how to learn such a *binary* mixture source from its $(2k - 1)$ -snapshots (see Algorithm 2 and Theorem 5.3). Thus, if we can obtain $(2k - 1)$ -snapshots from the binary source $(\vartheta, (f^1, \dots, f^k))$ (although our input is θ) then Theorem 5.3 would yield the desired result. We show that this is indeed possible, and hence, obtain the following result (whose proof appears at the end of this section).

Theorem 5.1. *Let $\theta = (\vartheta, (q^1, \dots, q^k))$ be a k -mixture source on $[0, 1]$, and (ϑ, α) be the corresponding k -spike distribution. Let $\tau = \min_{j \neq j'} |\alpha_j - \alpha_{j'}|$. For any $s < \tau$ and $\psi > 0$, using $3k2^{4k}s^{-4k} \ln(4k/\psi)$ $(2k - 1)$ -snapshots from source θ , one can compute in polytime a k -spike distribution $(\tilde{\vartheta}, \tilde{\alpha})$ on $[0, 1]$ such that $\text{Tran}(\vartheta, \alpha; \tilde{\vartheta}, \tilde{\alpha}) \leq 1024ks^{1/(4k)}$ with probability at least $1 - \psi$.*

5.1 Learning a binary k -mixture source

Recall that $(\vartheta, (f^1, \dots, f^k))$ denotes the binary k -mixture source, and $\alpha_i = f_1^i$ is the bias of the i -th “coin”. We can collect from each $(2k - 1)$ -snapshot a random variable $0 \leq X \leq 2k - 1$ denoting the number of times the outcome “1” occurs in the snapshot. Thus,

$$\Pr[X = i] = \binom{2k-1}{i} \sum_{j=1}^k \vartheta_j \alpha_j^i (1 - \alpha_j)^{2k-1-i}. \quad (4)$$

Our objective is to use these statistics to reconstruct, in transportation distance (see Section 2.2), the binary source (i.e., the mixture weights and the k biases). Now consider the equivalent k -spike distribution (ϑ, α) . The i -th moment, and (what we call) the i -th *normalized binomial*

moment (NBM) of this distribution are

$$g_i(\vartheta, \alpha) = \sum_{j=1}^k \vartheta_j \alpha_j^i, \quad \nu_i(\vartheta, \alpha) = \sum_{j=1}^k \vartheta_j \alpha_j^i (1 - \alpha_j)^{2k-1-i}$$

Up to the factors $\binom{2k-1}{i}$ the NBMs are precisely the statistics of the random variable X (Eqn. 4) and so our objective in this section can be restated as: use the empirical NBM's to reconstruct the k -spike distribution (ϑ, α) .

Let $g(\vartheta, \alpha) = (g_i(\vartheta, \alpha))_{i=0}^{2k-1}$ and $\nu(\vartheta, \alpha) = (\nu_i(\vartheta, \alpha))_{i=0}^{2k-1}$ denote the vector of the first $2k-1$ moments and NBMs respectively of (ϑ, α) . For a positive integer b , and a vector $\beta = (\beta_1, \dots, \beta_\ell)$, let $A_b(\beta)$ be the $\ell \times b$ matrix $(A_b(\beta))_{ij} = (1 - \beta_i)^{b-1-j} \beta_i^j$ (with $1 \leq i \leq \ell$ and $0 \leq j \leq b-1$). Analogously, let $V_b(\beta)$ be the $\ell \times b$ matrix $(V_b(\beta))_{ij} = \beta_i^j$ (with $1 \leq i \leq \ell$ and $0 \leq j \leq b-1$). Let Pas be the $2k \times 2k$ lower triangular ‘‘Pascal triangle’’ matrix: for $0 \leq j \leq 2k-1$ and $j+1 \leq i \leq 2k$, $\text{Pas}_{ij} = \binom{2k-j-1}{i-j-1}$. Then $V_{2k}(\alpha) = A_{2k}(\alpha)\text{Pas}$, $\nu(\vartheta, \alpha) = \vartheta A_{2k}(\alpha)$, and $g(\vartheta, \alpha) = \vartheta V_{2k}(\alpha) = \nu(\vartheta, \alpha)\text{Pas}$.

In our algorithm it is convenient to use the empirical ordinary moments, but what we obtain are actually the empirical NBM's, so we need the following lemma.

Lemma 5.2. $\|\text{Pas}\|_{\text{op}} \leq 4^k / \sqrt{3}$.

Proof. $\|\text{Pas}\|_{\text{op}} \leq \|\text{Pas}\|_F = \sqrt{\sum_{m=0}^{2k-1} \sum_{i=0}^m \binom{m}{i}^2}$. Since $\sum_{i=0}^m \binom{m}{i}^2 = \binom{2m}{m} \leq 2^{2m}$, $\|\text{Pas}\|_F \leq \sqrt{\sum_{m=0}^{2k-1} 2^{2m}}$. ■

Our algorithm uses two input parameters τ and ξ as input, and the empirical NBM vector $\tilde{\nu}$ (or equivalently \tilde{g}). Since we infer (in the sampling limit) the locations of the k spikes exactly, there is a singularity in the process when spikes coincide. So we assume a minimum separation between spikes: $\tau = \min_{j \neq j'} |\alpha_j - \alpha_{j'}|$. (It is of course possible to simply run a doubling search for sufficiently small τ , but the required accuracy in the moments, and hence sample size, does increase as τ decreases.) We also assume a bound ξ on the accuracy of our empirical statistics. (When we utilize Theorem 5.3 to obtain Theorem 5.1, ξ is a consequence, and not an input parameter). We require that

$$\|\tilde{\nu} - \nu(\vartheta, \alpha)\|_2 \leq \xi 4^{-k} \sqrt{3}, \quad \xi \leq \tau^{2k} \tag{5}$$

Theorem 5.3. *There is a polytime algorithm that receives as input τ, ξ , an empirical NBM vector $\tilde{\nu} \in \mathbb{R}^{2k}$ satisfying (5), and outputs a k -spike distribution $(\tilde{\vartheta}, \tilde{\alpha})$ on $[0, 1]$ such that $\text{Tran}(\vartheta, \alpha; \tilde{\vartheta}, \tilde{\alpha}) \leq O(\xi^{\Omega(1/k^2)})$.*

We first show the information-theoretic feasibility of Theorem 5.3: the transportation distance between two probability measures on $[0, 1]$ is upper bounded by (a moderately-growing function of) the Euclidean distance between their moment maps. (To use Lemma 5.4 to prove Theorem 5.3, we have to show how to compute $\tilde{\vartheta}$ and $\tilde{\alpha}$ from \tilde{g} such that $\|\tilde{g} - g(\tilde{\vartheta}, \tilde{\alpha})\|_2$, and hence, $\|g(\vartheta, \alpha) - g(\tilde{\vartheta}, \tilde{\alpha})\|_2$ is small.)

Lemma 5.4. *For any two (at most) k -spike distributions (ϑ, α) $(\tilde{\vartheta}, \tilde{\alpha})$ on $[0, 1]$,*

$$\|g(\vartheta, \alpha) - g(\tilde{\vartheta}, \tilde{\alpha})\|_2 \geq \frac{1}{(2k-1)^{4k} 2^{8k-5}} \cdot \left(\text{Tran}(\vartheta, \alpha; \tilde{\vartheta}, \tilde{\alpha}) \right)^{4k-2}.$$

Lemma 5.4 can be geometrically interpreted as follows. The point $g(\vartheta, \alpha)$ is in the convex hull of the moment curve and is therefore, by Caratheodory's theorem, expressible as a convex combination

of $2k$ points on the curve. However, this point is special in that it belongs to the collection of points expressible as a convex combination of merely k points of the curve. Lemma 5.4 shows that $g(\vartheta, \alpha)$ is in fact *uniquely* expressible in this way, and that moreover this combination is stable: any nearby point in this collection can only be expressed as a very similar convex combination. We utilize the following lemma, which can be understood as a global curvature property of the moment curve; we defer its proof to Section 5.2. The moment curve plays a central role in convex and polyhedral geometry [8], but as far as we know Lemmas 5.4 and 5.5 are new, and may be of independent interest.

Lemma 5.5. *Let $0 \leq \beta_1 < \dots < \beta_{\kappa+1} \leq 1$, $\ell \in [\kappa]$, and $s = \beta_{\ell+1} - \beta_\ell$. Let $\gamma(x) = \sum_{i=0}^\kappa \gamma_i x^i$ be a real polynomial of degree κ evaluating to 1 at the points $\beta_1, \dots, \beta_\ell$ and evaluating to 0 at the points $\beta_{\ell+1}, \dots, \beta_{\kappa+1}$. Then $\sum_{i=0}^\kappa \gamma_i^2 \leq \kappa^2 2^{4\kappa-1} s^{-2\kappa}$.*

Proof of Lemma 5.4. Denote $\{\alpha_1, \dots, \alpha_k\} \cup \{\tilde{\alpha}_1, \dots, \tilde{\alpha}_k\}$ by $\bar{\alpha} = \{\bar{\alpha}_1, \dots, \bar{\alpha}_K\}$ where $\bar{\alpha}_1 < \dots < \bar{\alpha}_K$. Define $\bar{\vartheta}_i = \sum_{j:\alpha_j=\bar{\alpha}_i} \vartheta_j - \sum_{j:\tilde{\alpha}_j=\bar{\alpha}_i} \tilde{\vartheta}_j$ for $i \in [K]$. Let $\bar{\vartheta} \in \mathbb{R}^K$ be the row vector $(\bar{\vartheta}_1, \dots, \bar{\vartheta}_K)$. Let $\eta = \text{Tran}(\vartheta, \alpha; \tilde{\vartheta}, \tilde{\alpha})$. So we need to show that $\|\bar{\vartheta} V_{2k}(\bar{\alpha})\|_2 \geq \frac{1}{(2k-1)^{4k} 2^{8k-5}} \cdot \eta^{4k-2}$. It suffices to show that $\|\bar{\vartheta} V_K(\bar{\alpha})\|_2 \geq \frac{1}{(K-1)^{2K} 2^{4K-5}} \cdot \eta^{2K-2}$.

There is an $1 \leq \ell < K$ such that $|\sum_{i=1}^\ell \bar{\vartheta}_i| \cdot (\bar{\alpha}_{\ell+1} - \bar{\alpha}_\ell) \geq \eta/(K-1)$. Let $\delta = \sum_{i=1}^\ell \bar{\vartheta}_i$; without loss of generality $\delta \geq 0$, and note that $\delta \leq 1$. Let $s = \bar{\alpha}_{\ell+1} - \bar{\alpha}_\ell$, so $(K-1)\delta s \geq \eta$.

Denote row i of a matrix Z by Z_{i*} and column j by Z_{*j} . We lower bound $\|\bar{\vartheta} V_K(\bar{\alpha})\|_2$, by considering its minimum value under the constraints $\sum_{i=1}^\ell \bar{\vartheta}_i = \delta$ and $\sum_{i=1}^K \bar{\vartheta}_i = 0$.

A vector $y^\dagger = \bar{\vartheta} V_K(\bar{\alpha})$ minimizing $\|y\|_2$ must be orthogonal to $V_K(\bar{\alpha})_{i*} - V_K(\bar{\alpha})_{i' *}$ if $1 \leq i < i' \leq \ell$ or if $\ell+1 \leq i < i' \leq K$. This means that there are scalars c and d such that $V_K(\bar{\alpha})y = c(\sum_{j=1}^\ell e_j) + d(\sum_{j=\ell+1}^K e_j)$, where vector $e_j \in \mathbb{R}^K$ has a 1 in the j -th position and 0 everywhere else. Therefore, $y = c\gamma + d\gamma'$, where $\gamma = \sum_{j=1}^\ell (V_K(\bar{\alpha})^{-1})_{*j}$ and $\gamma' = \sum_{j=\ell+1}^K (V_K(\bar{\alpha})^{-1})_{*j}$. At the same time

$$\delta = \sum_{i=1}^\ell \bar{\vartheta}_i = \bar{\vartheta} V_K(\bar{\alpha})\gamma = y^\dagger \gamma = c\|\gamma\|_2^2 + d\gamma^\dagger \gamma \quad -\delta = \sum_{i=\ell+1}^K \bar{\vartheta}_i = \bar{\vartheta} V_K(\bar{\alpha})\gamma' = y^\dagger \gamma' = c\gamma^\dagger \gamma' + d\|\gamma'\|_2^2$$

and hence, $\|y\|_2^2 = y \cdot (c\gamma + d\gamma') = (c-d)\delta$. Solving for c, d ,

$$c-d = \frac{\delta\|\gamma + \gamma'\|_2^2}{\|\gamma\|_2^2 \cdot \|\gamma'\|_2^2 - (\gamma^\dagger \cdot \gamma')^2}.$$

First we examine the numerator. Like any combination of the columns of $V_K(\bar{\alpha})^{-1}$, $\gamma + \gamma'$ is the list of coefficients of a polynomial of degree $K-1$, in the basis $1, x, \dots, x^{K-1}$. By definition, $\gamma + \gamma' = \sum_j (V_K(\bar{\alpha})^{-1})_{*j}$, which is to say that for every i , $V_K(\bar{\alpha})_{i*} \cdot (\gamma + \gamma') = 1$. So the polynomial $\gamma + \gamma'$ evaluates to 1 at every $\bar{\alpha}_i$. It can therefore only be the constant polynomial 1; this means that $(\gamma + \gamma')_i = 1$ if $i = 0$, and $(\gamma + \gamma')_i = 0$ otherwise. Thus $\|\gamma + \gamma'\|_2^2 = 1$.

Next we examine the denominator, which we upper bound by $\|\gamma\|_2^2 \cdot \|\gamma'\|_2^2$. When interpreted as a polynomial, γ takes the value 1 on a nonempty set of points $\bar{\alpha}_1, \dots, \bar{\alpha}_\ell$ separated by the positive distance $s = \bar{\alpha}_{\ell+1} - \bar{\alpha}_\ell$ from another nonempty set of points $\bar{\alpha}_{\ell+1}, \dots, \bar{\alpha}_K$ upon which it takes the value 0. Observe that if the polynomial was required to change value by a large amount within a short interval, it would have to have large coefficients. A converse to this is the inequality stated in Lemma 5.5. Using this to bound $\|\gamma\|_2^2$ and $\|\gamma'\|_2^2$, and since $\delta s \geq \eta/(K-1)$, we obtain that

$$\|y\|_2^2 = (c-d)\delta \geq \frac{\delta^2}{\|\gamma\|_2^2 \cdot \|\gamma'\|_2^2} \geq \frac{\delta^2}{((K-1)^{2K} 2^{4K-5} s^{-2K+2})^2} \geq \frac{\eta^{4K-4}}{(K-1)^{4K} 2^{8K-10}}. \quad \blacksquare$$

We now define the algorithm promised by Theorem 5.3, and then prove the theorem. To give some intuition, suppose first that we are given the true moment vector $g(\vartheta, \alpha) = \vartheta V_{2k}(\alpha)$. Observe that there is a common vector $\lambda = (\lambda_0, \dots, \lambda_k)^\top$ of length $k+1$ that is a dependency among every $k+1$ adjacent columns of $V_{2k}(\alpha)$. In other words, letting $\Lambda = \Lambda(\lambda)$ denote the $2k \times k$ matrix with $\Lambda_{ij} = \lambda_{i-j}$ (for $0 \leq i < 2k$, $0 \leq j < k$ and with the understanding $\lambda_\ell = 0$ for $\ell \notin \{0, \dots, k\}$), $V_{2k}(\alpha)\Lambda = 0$. Thus $g(\vartheta, \alpha)\Lambda = \vartheta V_{2k}(\alpha)\Lambda = 0$. Overtly this is a system of $2k$ equations to determine λ . But we eliminate the redundancy in Λ by forming the $k \times (k+1)$ matrix $G = G(g(\vartheta, \alpha))$ defined by $G_{ij} = g(\vartheta, \alpha)_{i+j}$ for $i = 0, \dots, k-1$ and $j = 0, \dots, k$, and then solve the system of linear equations $G\lambda = 0$ to obtain λ . This system does not have a unique solution, so in the sequel λ will denote a solution with $\lambda_k = 1$. For each $i = 1, \dots, k$, we have $(V_{2k}(\alpha)\Lambda)_{i,1} = \sum_{\ell=0}^k \lambda_\ell \alpha_i^\ell = 0$. This implies that we can obtain the α_i values by computing the roots of the polynomial $P_\lambda(x) := \sum_{\ell=0}^k \lambda_\ell x^\ell$. Once we have the α_i 's, we can compute ϑ by solving for y the system of linear equations $yV_{2k}(\alpha) = g(\vartheta, \alpha)$.

Of course, we are actually given \tilde{g} rather than the true vector $g(\vartheta, \alpha)$. So we need to control the error in estimating first α and then ϑ . The learning algorithm is as follows.

Algorithm 2. *Input: parameters ξ, τ and empirical moments \tilde{g} such that $\|tg - g(\vartheta, \alpha)\|_2 \leq \xi$.
Output: a k -spike distribution $(\tilde{\vartheta}, \tilde{\alpha})$*

B1. Solve the minimization problem:

$$\text{minimize } \|x\|_1 \quad \text{s.t.} \quad \|G(\tilde{g})x\|_1 \leq 2^k k \xi, \quad x_k = 1 \quad (\text{P})$$

which can be encoded as a linear program, to obtain a solution $\tilde{\lambda}$. Observe that since $G(\tilde{g})$ has $k+1$ columns and k rows, there is always a feasible solution.

B2. Let $\bar{\alpha}_1, \dots, \bar{\alpha}_k$ be the (possibly complex) roots of the polynomial $P_{\tilde{\lambda}}$. Thus, we have $V_{2k}(\bar{\alpha})\Lambda(\tilde{\lambda}) = 0$. We map the roots to values in $[0, 1]$ as follows. Let $\epsilon = \frac{4}{\tau}(2k\xi)^{1/k}$. First we compute values $\hat{\alpha}_1, \dots, \hat{\alpha}_k$ such that $|\hat{\alpha}_i - \bar{\alpha}_i| \leq \epsilon$ for every i , in time $\text{poly}(\log(\frac{1}{\epsilon}))$, using Pan's algorithm [35, Theorem 1.1]². We now set $\tilde{\alpha}_i = \max\{0, \min\{\text{Re}(\hat{\alpha}_i), 1\}\}$.

B3. Finally, we set $\tilde{\vartheta}$ to be the row-vector y that minimizes $\|yV_{2k}(\tilde{\alpha}) - \tilde{g}\|_2$ subject to $\|y\|_1 = 1, y \geq 0$. Note that this is a convex program.

We now analyze Algorithm 2 and justify Theorem 5.3. Recall that $\tau = \min_{j \neq j'} |\alpha_j - \alpha_{j'}|$. We need the following lemma, whose proof appears in Section 5.3.

Lemma 5.6. *The weights $\tilde{\vartheta}$ satisfy $\|\tilde{\vartheta}V_{2k}(\tilde{\alpha}) - \tilde{g}\|_2 \leq \|g(\vartheta, \alpha) - \tilde{g}\|_2 + \frac{(8k)^{5/2}}{\tau} \cdot (2k\xi)^{1/k}$.*

Proof of Theorem 5.3. We call Algorithm 2 with $\tilde{g} = \tilde{\nu}\text{Pas}$. By Lemma 5.2, we obtain that $\|\tilde{g} - g(\vartheta, \alpha)\|_2 \leq \xi$, and by Lemma 5.6, we have that $\|g(\vartheta, \alpha) - \tilde{\vartheta}V_{2k}(\tilde{\alpha})\|_2 \leq 2\|g(\vartheta, \alpha) - \tilde{g}\|_2 + \frac{8}{\tau} \cdot (8k)^{3/2}(2k\xi)^{1/k}$. Coupled with Lemma 5.4 and since $\xi \leq \tau^{2k}$, we obtain that

$$\begin{aligned} \text{Tran}(\vartheta, \alpha; \tilde{\vartheta}, \tilde{\alpha}) &\leq \left[(2k-1)^{4k} 2^{8k-5} \|g(\vartheta, \alpha) - g(\tilde{\vartheta}, \tilde{\alpha})\|_2 \right]^{\frac{1}{4k-2}} \\ &\leq \left[(2k-1)^{4k} 2^{8k-5} \left(2\xi + \frac{(8k)^{5/2}}{\tau} (2k\xi)^{1/k} \right) \right]^{\frac{1}{4k-2}} \\ &\leq \left[(2k-1)^{4k} 2^{8k-5} \left(2\xi + (8k)^{5/2} (2k\sqrt{\xi})^{1/2k} \right) \right]^{\frac{1}{4k-2}} \\ &\leq 1024 \cdot k\xi^{\frac{1}{8k^2}}. \end{aligned} \quad \blacksquare$$

²The theorem requires that the complex roots lie within the unit circle and that the coefficient of the highest-degree term is 1; but the discussion following it in [35] shows that this is essentially without loss of generality.

Proof of Theorem 5.1. We convert θ to the corresponding binary source $(\vartheta, (f^1, \dots, f^k))$ by randomized rounding. Given a $(2k-1)$ -snapshot $z = (z_1, \dots, z_{2k-1}) \in [0, 1]^{2k-1}$ from θ , we obtain a $(2k-1)$ -snapshot from the binary source as follows. We choose $2k-1$ independent values a_1, \dots, a_{2k-1} uniformly at random from $[0, 1]$ and set $X_i = 1$ if $z_i \geq a_i$ and 0 otherwise for all $i \in [2k-1]$. Note that if q^j is the constituent generating the $(2k-1)$ -snapshot z , then $\Pr[X_i = 1|q^j] = \mathbb{E}[X_i|q^j] = \alpha_j$, and so X_1, \dots, X_{2k-1} is a random $(2k-1)$ -snapshot from the above binary source.

Now we apply Theorem 5.3, setting $\xi = s^{2k}$. Let $\tilde{\nu}$ be the empirical NBM-vector obtained from the $(2k-1)$ -snapshots of the above binary source (i.e., $\tilde{\nu}_i = \binom{2k-1}{i}^{-1}$ · (frequency with which the $(2k-1)$ -snapshot has exactly i 1s)). The stated sample size ensures, via a Chernoff bound, that $\Pr[|\tilde{\nu}_i - \nu(\vartheta, \alpha)_i| \geq \frac{\xi 4^{-k}}{\sqrt{6k}}] < \frac{\psi}{2k}$ for all $i = 0, \dots, 2k-1$. Hence, with probability at least $1 - \psi$, we have $\|\tilde{\nu} - \nu(\vartheta, \alpha)\|_2 \leq \sqrt{2k} \cdot \|\tilde{\nu} - \nu(\vartheta, \alpha)\|_\infty \leq \xi 4^{-k} / \sqrt{3}$. ■

5.2 Proof of Lemma 5.5

There are two easy cases to dismiss before we reach the more subtle part of this lemma. The first easy case is $\ell = 1$. In this case γ is a single Lagrange interpolant: $\gamma(x) = \prod_{j=2}^{\kappa+1} \frac{x - \beta_j}{\beta_1 - \beta_j}$. For $0 \leq i \leq \kappa$ let $e_i^\kappa(\beta_2, \dots, \beta_{\kappa+1})$ be the i 'th elementary symmetric mean,

$$e_i^\kappa(\beta_2, \dots, \beta_{\kappa+1}) = \frac{1}{\binom{\kappa}{i}} \sum_{S \subseteq \{2, \dots, \kappa+1\}: |S|=i} \prod_{j \in S} \beta_j$$

and observe that for all i , $0 \leq e_i^\kappa(\beta_2, \dots, \beta_{\kappa+1}) \leq 1$. Now

$$\gamma(x) = \left(\prod_{j=2}^{\kappa+1} \frac{1}{\beta_1 - \beta_j} \right) \sum_{i=0}^{\kappa} (-1)^{\kappa-i} \binom{\kappa}{i} e_{\kappa-i}^\kappa(\beta_2, \dots, \beta_{\kappa+1}) x^i$$

So $\sum_{i=0}^{\kappa} \gamma_i^2 = \left(\prod_{j=2}^{\kappa+1} \frac{1}{\beta_1 - \beta_j} \right)^2 \sum_{i=0}^{\kappa} \left(\binom{\kappa}{i} e_{\kappa-i}^\kappa(\beta_2, \dots, \beta_{\kappa+1}) \right)^2 \leq s^{-2\kappa} \sum_{i=0}^{\kappa} \binom{\kappa}{i}^2 = \binom{2\kappa}{\kappa} s^{-2\kappa}$.

The second easy case is $\ell = \kappa$; this is almost as simple. Merely note that the above argument applies to the polynomial $1 - \gamma$, so that we have only to allow for the possible increase of $|\gamma_0|$ by 1. Hence $\sum_{i=0}^{\kappa} \gamma_i^2 \leq 4 \binom{2\kappa}{\kappa} s^{-2\kappa}$.

We now consider the less trivial case of $1 < \ell < \kappa$. The difficulty here is that the Lagrange interpolants of γ may have very large coefficients, particularly if among $\beta_1, \dots, \beta_\ell$ or among $\beta_{\ell+1}, \dots, \beta_{\kappa+1}$ there are closely spaced roots, as well there may be. We must show that these large coefficients cancel out in γ .

The trick is to examine not γ but $\partial\gamma/\partial x$. The roots of the derivative interlace the two sets on which γ is constant, which is to say, with $\beta'_1 \leq \dots \leq \beta'_{\kappa-1}$ denoting the roots of $\partial\gamma/\partial x$, that for $j < \ell$, $\beta_j \leq \beta'_j \leq \beta_{j+1}$, and for $j \geq \ell$, $\beta_{j+1} \leq \beta'_j \leq \beta_{j+2}$. In particular, none of the roots fall in the interval $(\beta_\ell, \beta_{\ell+1})$. For some constant C we can write $\partial\gamma/\partial x = C \prod_{j=0}^{\kappa-1} (x - \beta'_j)$ (with $\text{sign}(C) = (-1)^{1+\kappa-\ell}$). Observe that $\int_{\beta_\ell}^{\beta_{\ell+1}} \frac{\partial\gamma}{\partial x}(x) dx = -1$. So $(-1)^{1+\kappa-\ell}/C = \int_{\beta_\ell}^{\beta_{\ell+1}} (-1)^{\kappa-\ell} \prod_{j=0}^{\kappa-1} (x - \beta'_j) dx$. Observe that if for any $j < \ell$, β'_j is increased, or if for any $j \geq \ell$, β'_j is decreased, then the integral decreases. So $(-1)^{1+\kappa-\ell}/C \geq \int_{\beta_\ell}^{\beta_{\ell+1}} (-1)^{\kappa-\ell} (x - \beta_\ell)^{\ell-1} (x - \beta_{\ell+1})^{\kappa-\ell} dx$. This is a definite integral that can be evaluated in closed form:

$$\int_{\beta_\ell}^{\beta_{\ell+1}} (-1)^{\kappa-\ell} (x - \beta_\ell)^{\ell-1} (x - \beta_{\ell+1})^{\kappa-\ell} dx = (\beta_{\ell+1} - \beta_\ell)^\kappa (\ell-1)! (\kappa-\ell)! / \kappa! .$$

Hence, $(-1)^{1+\kappa-\ell}C \leq \frac{\kappa!}{s^\kappa(\ell-1)!(\kappa-\ell)!}$. The sum of squares of coefficients of $\frac{\partial\gamma}{\partial x}$ is $C^2 \sum_{i=0}^{\kappa-1} \binom{\kappa-1}{i}^2 (e_i^{\kappa-1}(\beta'_1, \dots, \beta'_{\kappa-1}))^2 \leq C^2 \binom{2\kappa-2}{\kappa-1}$. Integration only decreases the magnitude of the coefficients, so the same bound applies to γ , with the exception of the constant coefficient. The constant coefficient can be bounded by the fact that γ has a root in $(0, 1)$, and that in that interval the derivative is bounded in magnitude by $C \sum_{i=0}^{\kappa-1} \binom{\kappa-1}{i} = C \cdot 2^\kappa$. So $|\gamma_0| \leq C \cdot 2^\kappa$. Consequently, $\sum_{i=0}^\kappa \gamma_i^2$ is at most

$$C^2 \left[\binom{2\kappa-2}{\kappa-1} + 2^{2\kappa} \right] \leq \frac{\binom{2\kappa-2}{\kappa-1} + 2^{2\kappa}}{s^{2\kappa}} \cdot \left(\frac{\kappa!}{(\ell-1)!(\kappa-\ell)!} \right)^2 \leq \frac{5\kappa^2 2^{2\kappa-2}}{s^{2\kappa}} \cdot \binom{\kappa-1}{\ell-1}^2 \leq \frac{5\kappa^2 2^{4\kappa-4}}{s^{2\kappa}},$$

which completes the proof of the lemma. \blacksquare

5.3 Proof of Lemma 5.6

Recall that $G = G(g(\vartheta, \alpha))$ is the $k \times (k+1)$ matrix defined by $G_{ij} = g(\vartheta, \alpha)_{i+j}$ for $i = 0, \dots, k-1$ and $j = 0, \dots, k$; λ is such that $G\lambda = 0$ and $\lambda_k = 1$; $\Lambda = \Lambda(\lambda)$ is the $2k \times k$ matrix with $\Lambda_{ij} = \lambda_{i-j}$ (for $0 \leq i < 2k$, $0 \leq j < k$ with the understanding $\lambda_\ell = 0$ for $\ell \notin \{0, \dots, k\}$; and $P_\lambda(x)$ is the polynomial $\sum_{\ell=0}^k \lambda_\ell x^\ell$. We use V_k, V_{2k} , to denote $V_k(\alpha), V_{2k}(\alpha)$ respectively, and $\tilde{V}_k, \tilde{V}_{2k}, \tilde{G}, \tilde{\Lambda}$ to denote $V_k(\tilde{\alpha}), V_{2k}(\tilde{\alpha}), G(\tilde{g}), \Lambda(\tilde{\lambda})$ respectively. We abbreviate $g(\vartheta, \alpha)$ to g .

Lemma 5.7. *If $\|\tilde{g} - g\|_2 \leq \xi$, then $\|G\tilde{\lambda}\|_1 \leq 2^{k+1}k\xi$.*

Proof. First, observe that $\tilde{G}\lambda = G\lambda + (\tilde{G} - G)\lambda = (\tilde{G} - G)\lambda$. Also $\|\lambda\|_2 \leq \|\lambda\|_1 = \prod_{i=1}^k (1 + \alpha_i) \leq 2^k$. The last two inequalities follows since $P_\lambda(x) = \prod_{i=1}^k (x - \alpha_i)$, and $P_\lambda(-1) = (-1)^k \|\lambda\|_1$. So for any $i = 1, \dots, k$, $|(G - \tilde{G})_i \cdot \lambda| \leq \|\lambda\|_2 \|G_i - \tilde{G}_i\|_2 \leq 2^k \xi$. Thus, λ is a feasible solution to (P), which implies that $\|\tilde{\lambda}\|_1 \leq 2^k$. We have $\|G\tilde{\lambda}\|_1 \leq \|\tilde{G}\tilde{\lambda}\|_1 + \|(G - \tilde{G})\tilde{\lambda}\|_1 \leq 2^k k\xi + \|(G - \tilde{G})\tilde{\lambda}\|_1$. For any $i = 1, \dots, k$, $|(G - \tilde{G})_i \cdot \tilde{\lambda}| \leq \|G_i - \tilde{G}_i\|_2 \|\tilde{\lambda}\|_2 \leq 2^k \xi$, so $\|G\tilde{\lambda}\|_1 \leq 2^{k+1}k\xi$. \blacksquare

Lemma 5.8. *For every α_i , $i = 1, \dots, k$, there exists a $\sigma(i) \in \{1, \dots, k\}$ such that $\vartheta_i |\alpha_i - \tilde{\alpha}_{\sigma(i)}| \leq \frac{8}{\tau} (2k\xi)^{1/k}$.*

Proof. Since $\|G\tilde{\lambda}\|_2 \leq 2^{k+1}k\xi$ (by Proposition 5.7), we have equivalently that the $\|\cdot\|_2$ norm of $g\tilde{\Lambda} = \vartheta V_{2k} \tilde{\Lambda}$ is at most $2^{k+1}k\xi$. We may write $\vartheta V_{2k} \tilde{\Lambda}$ as

$$\vartheta V_{2k} \tilde{\Lambda} = \vartheta \begin{pmatrix} P_{\tilde{\lambda}}(\alpha_1) & \alpha_1 P_{\tilde{\lambda}}(\alpha_1) & \cdots & \alpha_1^{k-1} P_{\tilde{\lambda}}(\alpha_1) \\ P_{\tilde{\lambda}}(\alpha_2) & \alpha_2 P_{\tilde{\lambda}}(\alpha_2) & \cdots & \alpha_2^{k-1} P_{\tilde{\lambda}}(\alpha_2) \\ \vdots & \vdots & \ddots & \vdots \\ P_{\tilde{\lambda}}(\alpha_k) & \alpha_k P_{\tilde{\lambda}}(\alpha_k) & \cdots & \alpha_k^{k-1} P_{\tilde{\lambda}}(\alpha_k) \end{pmatrix}$$

which is equal to $\vartheta' V_k(\alpha)$ where $\vartheta' = (\vartheta_1 P_{\tilde{\lambda}}(\alpha_1), \dots, \vartheta_k P_{\tilde{\lambda}}(\alpha_k))$. Thus, we are given that $\|\vartheta' V_k\|_2 \leq 2^{k+1}k\xi$.

Let $(\gamma^i)^\dagger = (\arg \min_{y \in \mathbb{R}^k: y_i=1} \|y V_k\|_2) V_k$. Then, we also have $\|\vartheta' V_k\|_2 \geq \max_i |\vartheta'_i| \|\gamma^i\|_2$. Note that γ^i must be orthogonal to $(V_k)_{j*}$ for all $j \neq i$, and $(V_k)_{i*} \gamma^i = \|\gamma^i\|_2^2$. (Recall that Z_{i*} denotes row i of a matrix Z .) Let $Q_i(x) = \sum_{\ell=0}^{k-1} \gamma_\ell^i x^\ell$. Then, $Q_i(x) = \|\gamma^i\|_2^2 \prod_{j \neq i} \frac{x - \alpha_j}{\alpha_i - \alpha_j}$. Also, since the coefficients of $Q_i(x)$ have alternating signs, we have

$$|Q_i(-1)| = \|\gamma^i\|_1 = \|\gamma^i\|_2^2 \prod_{j \neq i} \frac{1 + \alpha_j}{|\alpha_i - \alpha_j|}.$$

Hence, $\|\gamma^i\|_2 \geq \prod_{j \neq i} \frac{|\alpha_i - \alpha_j|}{1 + \alpha_j}$. So we obtain the lower bound

$$\|\vartheta' V_k\|_2 \geq \max_i \left(|\vartheta'_i| \cdot \prod_{j \neq i} \frac{|\alpha_i - \alpha_j|}{1 + \alpha_j} \right) \geq \max_i \left(\vartheta_i \left(\frac{\tau}{2} \right)^{k-1} \prod_{j=1}^k |\alpha_i - \bar{\alpha}_j| \right) \geq \max_i \left(\vartheta_i \left(\frac{\tau}{2} \right)^{k-1} \prod_{j=1}^k |\alpha_i - \operatorname{Re}(\bar{\alpha}_j)| \right).$$

The last inequality follows since complex roots occur in conjugate pairs, so if $\bar{\alpha}_\ell = a + bi$ is complex, then there must be some ℓ' such that $\bar{\alpha}_{\ell'} = a - bi$ and therefore,

$$\prod_j |\alpha_i - \bar{\alpha}_j| = ((\alpha_i - a)^2 + b^2) \cdot \prod_{j \neq \ell, \ell'} |\alpha_i - \bar{\alpha}_j| \geq (\alpha_i - a)^2 \cdot \prod_{j \neq \ell, \ell'} |\alpha_i - \bar{\alpha}_j|.$$

Now, we claim that $|\alpha_i - \operatorname{Re}(\bar{\alpha}_j)| \geq |\alpha_i - \tilde{\alpha}_j| - \epsilon$ for every j . If both $\operatorname{Re}(\bar{\alpha}_j)$ and $\operatorname{Re}(\hat{\alpha}_j)$ lie in $[0, 1]$, or both of them are less than 0, or both are greater than 1, then this follows since $|\bar{\alpha}_j - \hat{\alpha}_j| \leq \epsilon$ and $\alpha_i \in [0, 1]$. If $\operatorname{Re}(\bar{\alpha}_j) \notin [0, 1]$ but $\operatorname{Re}(\hat{\alpha}_j) \in [0, 1]$, or if $\operatorname{Re}(\bar{\alpha}_j) \in [0, 1]$ but $\operatorname{Re}(\hat{\alpha}_j) \notin [0, 1]$, then this again follows since $|\bar{\alpha}_j - \hat{\alpha}_j| \leq \epsilon$. Combining everything, we get that

$$2^k (2k\xi) \geq \|\vartheta' V_k\|_2 \geq \max_i \left(\vartheta_i \left(\frac{\tau}{2} \right)^{k-1} \prod_{j=1}^k |\alpha_i - \tilde{\alpha}_j| - \epsilon \right).$$

This implies that for every $i = 1, \dots, k$, there exists $\sigma(i) \in \{1, \dots, k\}$ such that $\vartheta_i |\alpha_i - \tilde{\alpha}_{\sigma(i)}| \leq \frac{4}{\tau} \cdot (2k\xi)^{1/k} + \epsilon$. ■

We can now wrap up the proof of Lemma 5.6. Let $\eta = \frac{8}{\tau} \cdot (2k\xi)^{1/k}$. We will bound $\|\tilde{\vartheta} \tilde{V}_{2k} - \tilde{g}\|_2$ by exhibiting a solution $y \in [0, 1]^k$, $\|y\|_1 = 1$ such that $\|y \tilde{V}_{2k} - \tilde{g}\|_2 \leq \|g - \tilde{g}\| + k(8k)^{3/2} \eta$. Let σ be the function whose existence is proved in Lemma 5.8. For $j = 1, \dots, k$, set $y_j = \sum_{i: \sigma(i)=j} \vartheta_i$ (if $\sigma^{-1}(j) = \emptyset$, then $y_j = 0$). We have $\|y \tilde{V}_{2k} - \tilde{g}\|_2 \leq \|g - \tilde{g}\|_2 + \|g - y \tilde{V}_{2k}\|_2$. We expand $g - y \tilde{V}_{2k} = \vartheta V_{2k} - y \tilde{V}_{2k} = \sum_{i=1}^k \vartheta_i ((V_{2k})_{i*} - (\tilde{V}_{2k})_{\sigma(i)*})$. For every i ,

$$\vartheta_i^2 \|(V_{2k})_{i*} - (\tilde{V}_{2k})_{\sigma(i)*}\|_2^2 = \vartheta_i^2 \sum_{\ell=0}^{2k-1} (\alpha_i^\ell - \tilde{\alpha}_{\sigma(i)}^\ell)^2 \leq \vartheta_i^2 \cdot 8k^3 \cdot \eta^2.$$

Therefore, $\|g - y \tilde{V}_{2k}\|_2 \leq k(8k)^{3/2} \eta$. ■

References

- [1] D. Achlioptas and F. McSherry. On spectral learning of mixtures of distributions. In *Proc. 18th Ann. Conf. on Learning Theory*, pages 458–469, June 2005.
- [2] N. Alon and J. H. Spencer. *The probabilistic method*. Wiley, 2nd edition, 2000.
- [3] A. Anandkumar, D.P. Foster, D. Hsu, S.M. Kakade, and Y.-K. Liu. Two SVDs suffice: Spectral decompositions for probabilistic topic modeling and latent Dirichlet allocation. *CoRR*, abs/1204.6703, 2012. <http://arxiv.org/abs/1204.6703>.
- [4] A. Anandkumar, D. Hsu, and S.M. Kakade. A method of moments for mixture models and hidden Markov models. In *Proc. 25th COLT*, pages 33.1–33.34, 2012. <http://arxiv.org/abs/1203.0683>.

- [5] S. Arora, R. Ge, and A. Moitra. Learning topic models — going beyond SVD. In *Proc. 53rd FOCS*, 2012. <http://arxiv.org/abs/1204.1956>.
- [6] S. Arora and R. Kannan. Learning mixtures of separated nonspherical Gaussians. *Ann. Appl. Prob.*, 15:69–92, 2005.
- [7] Y. Azar, A. Fiat, A. R. Karlin, F. McSherry, and J. Saia. Spectral analysis of data. In *Proc. 33rd STOC*, pages 619–626, 2001.
- [8] A. Barvinok. *A Course in Convexity*, volume 54 of *Graduate Studies in Mathematics*. AMS, 2002.
- [9] T. Batu, S. Guha, and S. Kannan. Inferring mixtures of Markov chains. In *Proc. 17th COLT*, pages 186–199, 2004.
- [10] M. Belkin and K. Sinha. Polynomial learning of distribution families. In *Proc. 51st FOCS*, pages 103–112, 2010. <http://doi.ieeecomputersociety.org/10.1109/FOCS.2010.16>.
- [11] D.M. Blei. Probabilistic topic models. *Commun. ACM*, 55(4):77–84, 2012. <http://doi.acm.org/10.1145/2133806.2133826>.
- [12] D.M. Blei, A.Y. Ng, and M.I. Jordan. Latent Dirichlet allocation. *J. Machine Learning Res.*, 3:993–1022, 2003. <http://www.jmlr.org/papers/v3/blei03a.html>.
- [13] S.C. Brubaker and S. Vempala. Isotropic PCA and affine-invariant clustering. In *FOCS*, pages 551–560, 2008. <http://doi.ieeecomputersociety.org/10.1109/FOCS.2008.48>.
- [14] K. Chaudhuri, E. Halperin, S. Rao, and S. Zhou. A rigorous analysis of population stratification with limited data. In *Proc. of the 18th SODA*, pages 1046–1055, 2007.
- [15] K. Chaudhuri and S. Rao. Beyond Gaussians: Spectral methods for learning mixtures of heavy-tailed distributions. In *Proc. 21st COLT*, pages 21–32, 2008. <http://colt2008.cs.helsinki.fi/papers/44-Chaudhuri.pdf>.
- [16] K. Chaudhuri and S. Rao. Learning mixtures of product distributions using correlations and independence. In *Proc. 21st COLT*, pages 9–20, 2008. <http://colt2008.cs.helsinki.fi/papers/7-Chaudhuri.pdf>.
- [17] M. Cryan, L. Goldberg, and P. Goldberg. Evolutionary trees can be learned in polynomial time in the two state general Markov model. *SICOMP*, 31(2):375–397, 2002.
- [18] A. Dasgupta, J. Hopcroft, J. Kleinberg, and M. Sandler. On learning mixtures of heavy-tailed distributions. In *Proc. 46th FOCS*, pages 491–500, 2005.
- [19] S. Dasgupta. Learning mixtures of Gaussians. In *Proc. of the 40th FOCS*, pages 634–644, 1999.
- [20] S. Dasgupta and L.J. Schulman. A probabilistic analysis of EM for mixtures of separated, spherical Gaussians. *Journal of Machine Learning Research*, 8:203–226, 2007. <http://www.jmlr.org/papers/v8/dasgupta07a.html>.
- [21] C. Daskalakis, I. Diakonikolas, and R.A. Servedio. Learning k -modal distributions via testing. In *Proc. 23rd SODA*, pages 1371–1385, 2012.

- [22] J. Feldman, R. O'Donnell, and R.A. Servedio. PAC learning mixtures of axis-aligned Gaussians with no separation assumption. In *Proc. 19th COLT*, pages 20–34, 2006.
- [23] J. Feldman, R. O'Donnell, and R.A. Servedio. Learning mixtures of product distributions over discrete domains. *SIAM J. Comput.*, 37(5):1536–1564, 2008. <http://dx.doi.org/10.1137/060670705>.
- [24] Y. Freund and Y. Mansour. Estimating a mixture of two product distributions. In *Proc. 12th Ann. Conf. on Computational Learning Theory*, pages 183–192, July 1999.
- [25] T. Hofmann. Probabilistic latent semantic analysis. In *Proc. 15th UAI*, pages 289–296, 1999. <http://uai.sis.pitt.edu/papers/99/p289-hofmann.pdf>.
- [26] T. Hofmann and J. Puzicha. Latent class models for collaborative filtering. In *Proc. IJCAI*, pages 688–693, 1999.
- [27] R. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 1985.
- [28] A.T. Kalai, A. Moitra, and G. Valiant. Efficiently learning mixtures of two Gaussians. In *Proc. 42nd STOC*, pages 553–562, June 2010.
- [29] R. Kannan, H. Salmasian, and S. Vempala. The spectral method for general mixture models. *SIAM J. Computing*, 38(3):1141–1156, 2008. <http://dx.doi.org/10.1137/S0097539704445925>.
- [30] M. Kearns, Y. Mansour, D. Ron, R. Rubinfeld, R. Schapire, and L. Sellie. On the learnability of discrete distributions. In *Proc. of the 26th STOC*, pages 273–282, 1994.
- [31] J. Kleinberg and M. Sandler. Using mixture models for collaborative filtering. *J. Comput. Syst. Sci.*, 74:49–69, 2008.
- [32] F. McSherry. Spectral partitioning of random graphs. In *Proc. 42nd FOCS*, pages 529–537, October 2001.
- [33] A. Moitra and G. Valiant. Settling the polynomial learnability of mixtures of Gaussians. In *Proc. of the 51st FOCS*, pages 93–102, 2010.
- [34] E. Mossel and S. Roch. Learning nonsingular phylogenies and hidden Markov models. In *STOC*, pages 366–375, 2005.
- [35] V. Y. Pan. Optimal and nearly optimal algorithms for approximating polynomial zeros. *Computers & Mathematics with Applications*, 31(12):97–138, 1996.
- [36] C. Papadimitriou, P. Raghavan, H. Tamaki, and S. Vempala. Latent semantic indexing: A probabilistic analysis. *J. Comput. Syst. Sci.*, 61(2):217–235, 2000.
- [37] Y. Rabani, L.J. Schulman, and C. Swamy. Inference from sparse sampling. <http://www.cs.technion.ac.il/~rabani/Papers/RabaniSS-manuscript.pdf>, 2008.
- [38] S. Vempala and G. Wang. A spectral algorithm for learning mixtures of distributions. *J. Comput. Syst. Sci.*, 68(4):841–860, 2004.
- [39] C. Villani. *Topics in optimal transportation*, volume 58 of *Graduate Studies in Mathematics*. AMS, 2003.
- [40] V. H. Vu. Spectral norm of random matrices. In *Proc. 37th STOC*, pages 423–430, 2005.

A Sample-size dependence of [3, 4, 5] on n for ℓ_1 -reconstruction

We view $P = (p^1, \dots, p^k)$ as an $n \times k$ matrix. Recall that $r = \sum_{t=1}^k w_t p_t p_t^\dagger$, $A = \sum_{t=1}^k w_t (p^t - r)(p^t - r)^\dagger$, and $M = rr^\dagger + A$. Let $w_{\max} := \max_t w_t$. We consider isotropic k -mixture sources, which is justified by Lemma 3.3. So $\frac{1}{2n} \leq r_i \leq \frac{2}{n}$ for all $i \in [n]$. Note that $\|r\|_1$ and $\|r\|_2^2$ are both $\Theta(\frac{1}{n})$. It will be convenient to split the width parameter ζ into two parameters. Let (i) $\frac{\zeta_1}{\sqrt{n}} = \min_{p, q \in P, p \neq q} \|p - q\|_2$; and (ii) $\frac{\zeta_2^2}{n}$ be the smallest non-zero eigenvalue of A . Then, the width of (w, P) is $\zeta = \max\{\zeta_1, \zeta_2\}$. We use $\sigma_i(Z)$ to denote the i -th largest singular value of a matrix Z . If Z has rank ℓ , its condition number is given by $\kappa(Z) := \sigma_1(Z)/\sigma_\ell(Z)$. For a square matrix Z with real eigenvalues, we use $\lambda_i(Z)$ to denote the i -th largest eigenvalue of Z . Note that if Z is an $n \times k$ matrix, then $\sigma_i(Z)^2 = \lambda_i(ZZ^\dagger) = \lambda_i(Z^\dagger Z)$ for all $i = 1, \dots, k$. Also the singular values of ZZ^\dagger coincide with its eigenvalues, and the same holds for $Z^\dagger Z$.

We now proceed to evaluate the sample-size dependence of [3, 4, 5] on n for reconstructing the mixture constituents within ℓ_1 -distance ϵ . Since these papers use different parameters than we do, in order to obtain a meaningful comparison, we relate their bounds to our parameters ζ_1, ζ_2 ; we keep track of the resulting dependence on n but ignore the (polynomial) dependence on other quantities. We show that the sample size needed is at least $\Omega(\frac{n^4}{\epsilon^2})$, with the exception of Algorithm B in [4], which needs $\Omega(\frac{n^3}{\epsilon^2})$ samples. As required by [3, 4, 5], we assume that P has full column rank. It follows that M has rank k and A has rank $k - 1$. The following inequality will be useful.

Proposition A.1. *Let $D = \text{diag}(d_1, \dots, d_k)$ where $d_1 \geq d_2 \geq \dots \geq d_k > 0$. Then $\lambda_k(PDP^\dagger) \geq d_k \lambda_k(PP^\dagger) = d_k \sigma_k(P)^2$.*

Comparison with [5]. The algorithm in [5] requires also that P be ρ -separable. This means that for every $t \in [k]$, there is some $i \in [n]$ such that $p_i^t \geq \rho$ and $p_i^{t'} = 0$ for all $t' \neq t$. This has the following implications. For any $t, t' \in [k]$, $t \neq t'$, we have $\|p^t - p^{t'}\|_2 \geq \sqrt{2}\rho$, so $\frac{\zeta_1}{\sqrt{n}} \geq \sqrt{2}\rho$. We can write $P^\dagger P = Y + Z$, where Y is a PSD matrix, and Z is a diagonal matrix whose diagonal entries are at least ρ^2 . So $\lambda_k(P^\dagger P) = \lambda_k(PP^\dagger) \geq \rho^2$. Therefore,

$$\frac{\zeta_2^2}{n} + \|r\|_2^2 = \lambda_k(A) + \|r\|_2^2 \geq \lambda_k(M) \geq w_{\min} \cdot \rho^2$$

where the first inequality follows from Lemma 2.2, and the second from Proposition A.1. It follows that $\rho = O(\frac{1}{\sqrt{n}})$. The bound in [5] to obtain ℓ_∞ error ϵ is (ignoring dependence on other quantities) $\Omega(\frac{1}{\epsilon^2 \rho^6})$. So setting $\epsilon = \frac{\epsilon}{n}$ to guarantee ℓ_1 -error at most ϵ and plugging in the above upper bounds on ρ , we obtain that the sample size is $\Omega(\frac{n^5}{\epsilon^2})$.

Comparison with [3]. The sample size required by [3] for the latent Dirichlet model for obtaining ℓ_2 error ϵ is $\Omega(\frac{1}{\epsilon^2 \sigma_k(P)^6})$. Proposition A.1 yields $\lambda_k(M) \geq w_{\min} \cdot \sigma_k(P)^2$ and as argued above, $\lambda_k(M) \leq \lambda_k(A) + \|r\|_2^2 = O(\frac{1}{n})$. So $\sigma_k(P)^6 = O(\frac{1}{n^3})$. Setting $\epsilon = \frac{\epsilon}{\sqrt{n}}$ for ℓ_1 error ϵ , this yields a bound of $\Omega(\frac{n^4}{\epsilon^2})$.

Comparison with [4]. Algorithm A in [4] requires sample size $\Omega(\frac{1}{\sigma_k(P)^8 \sigma_k(M)^4 \epsilon^2})$ to recover each p^t to within ℓ_2 -distance $\epsilon \max_{p \in P} \|p\|_2$. Since $\max_{p \in P} \|p\|_2 \leq \frac{2}{w_{\min} \sqrt{n}}$ due to isotropy, we can set $\epsilon = \frac{\epsilon w_{\min}}{2}$ to obtain ℓ_1 -error ϵ . Since $\sigma_k(P)^2$ and $\sigma_k(M) = \lambda_k(M)$ are both $O(\frac{1}{n})$, we obtain a bound of $\Omega(\frac{n^8}{\epsilon^2})$.

Algorithm B in [4] uses sample size $\Omega(\kappa(P)^8 / (\frac{\zeta_1^2}{n} \cdot \sigma_k(M)^2 \epsilon^2))$ to recover each p^t to within ℓ_2 -distance $\epsilon \max_{p \in P} \|p\|_2$. Clearly $\kappa(P) \geq 1$. Again, setting $\epsilon = \frac{\epsilon w_{\min}}{2}$, this yields a sample size of $\Omega(\frac{n^3}{\epsilon^2})$ for ℓ_1 error ϵ .